

# Conditional Gradient Algorithms for Rank-One Matrix Approximations with a Sparsity Constraint

Ronny Luss<sup>1</sup> and Marc Teboulle

School of Mathematical Sciences

Tel-Aviv University, Ramat-Aviv 69978, Israel

email: ronnyluss@gmail.com, teboulle@post.tau.ac.il

June 29, 2012

## Abstract

The sparsity constrained rank-one matrix approximation problem is a difficult mathematical optimization problem which arises in a wide array of useful applications in engineering, machine learning and statistics, and the design of algorithms for this problem has attracted intensive research activities. We introduce an algorithmic framework, called ConGradU, that unifies a variety of seemingly different algorithms that have been derived from disparate approaches, and allows for deriving new schemes. Building on the old and well-known conditional gradient algorithm, ConGradU is a simplified version with unit step size and yields a generic algorithm which either is given by an analytic formula or requires a very low computational complexity. Mathematical properties are systematically developed and numerical experiments are given.

**Keywords:** Sparse Principal Component Analysis, PCA, Conditional Gradient Algorithms, Sparse Eigenvalue Problems, Matrix Approximations

## 1 Introduction

The problem of interest here is the sparsity constrained rank-one matrix approximation given by

$$\max\{x^T A x : \|x\|_2 = 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\}, \quad (1)$$

where  $A \in \mathbf{S}^n$  is a given real symmetric matrix, and  $1 < k \leq n$  is a parameter controlling the sparsity of  $x$  which is defined by counting the number of nonzero entries of  $x$  and denoted using the  $l_0$  notation:  $\|x\|_0 = |\{i : x_i \neq 0\}|$ . This problem is also commonly known as the sparse Principal Component Analysis (PCA) problem, or as we refer to it,  $l_0$ -constrained PCA. Without the  $l_0$  constraint, the problem reduces to finding the first principal eigenvector and the corresponding maximal eigenvalue of the matrix  $A$ , i.e., solves

$$\max\{x^T A x : \|x\|_2 = 1, x \in \mathbf{R}^n\},$$

which is the PCA problem.

Suppose  $A = B^T B$  where  $B$  is an  $m \times n$  mean-centered data matrix with  $m$  samples and  $n$  variables, and denote by  $v$  be the principal eigenvector of  $A$ , i.e.,  $v$  solves the above PCA problem. Then  $Bv$  projects the data  $B$  to one dimension that maximizes the variance of the projected data. In general, PCA can be used to reduce  $B$  to  $l < n$  dimensions via the projection  $B(v_1, \dots, v_l)$  with  $v_i$  as the  $i^{th}$  eigenvector of  $A$ . In addition to dimensionality reduction, PCA can be used for visualization, clustering, and other tasks for

---

<sup>1</sup>Corresponding author

data analysis. Such tasks occur in various fields, e.g., genetics [1, 27], face recognition [16, 38], and signal processing [20, 17].

In PCA, the eigenvector is typically dense, i.e., each component of the eigenvector is nonzero, and hence the projected variables are linear functions of all original  $n$  variables. In sparse PCA, we restrict the number of variables used in this linear projection, thereby making it easier to interpret the projections. The additional  $l_0$  constraint however makes problem (1) a difficult and mathematically interesting problem which arises in many scientific and engineering applications where very large-scale data sets must be analyzed and interpreted.

Not surprisingly, the search and development of adequate algorithms for solving problem (1) have thus received much attention in the past decade, and this will be discussed below. But first, we want to make clear the main purpose of this paper. We have three main goals:

- To develop a novel and very simple approach to the  $l_0$ -constrained PCA problem (1) as formulated and without any modifications, i.e., no relaxations or penalization, which is amenable to dimensions in the hundreds of thousands or even millions.
- To present a “Father Algorithm”, which we call ConGradU, based on the well-known first-order conditional gradient scheme, which is very simple, allows for a rigorous convergence analysis, and provides a family of cheap algorithms well-suited to solving various formulations of sparse PCA.
- To provide a closure and unification to many seemingly disparate approaches recently proposed, and which will be shown to be particular realizations of ConGradU.

Most current approaches to sparse PCA can be categorized as solving one of several modified optimization problems based on penalization, relaxations, or both, and include:

- (a)  $l_1$ -constrained PCA:  $\max \{x^T A x : \|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}, x \in \mathbf{R}^n\}$ ,
- (b)  $l_0$ -penalized PCA:  $\max \{x^T A x - s\|x\|_0 : \|x\|_2 \leq 1, x \in \mathbf{R}^n\}$ ,
- (c)  $l_1$ -penalized PCA:  $\max \{x^T A x - s\|x\|_1 : \|x\|_2 \leq 1, x \in \mathbf{R}^n\}$ ,
- (d) Approximate  $l_0$ -penalized PCA:  $\max \{x^T A x - sg_p(x) : \|x\|_2 \leq 1, x \in \mathbf{R}^n\}$ ,  
where  $g_p(x) \simeq \|x\|_0$  and  $p$  controls the approximation,
- (e) Convex relaxations.

These models will be presented in detail in the forthcoming section, except for approach (e) which is not thoroughly discussed in this paper (but see §1.1 below).

In a nutshell, the original  $l_0$ -constrained PCA problem (1) and the corresponding modified problems (a)-(d) above can be written or transformed in such a way that they reduce to maximizing a convex function over some compact set  $C \subseteq \mathbf{R}^n$ :

$$(P) \quad \max\{F(x) : x \in C\}.$$

When the problem of maximizing a *linear* function over the compact set  $C$  can be efficiently computed, or even better, can be obtained *analytically*, a very simple and natural iterative scheme to consider for solving (P) is the so-called conditional gradient algorithm [21, 12]<sup>2</sup>.

---

<sup>2</sup>The conditional gradient scheme is also known as the Frank-Wolfe algorithm [15]. The latter was devised to minimize quadratic convex functions over a bounded polyhedron, while the former was extended mainly to solve convex minimization problems, see Section 3 for more precise details and relevant references.

To achieve the goals alluded to above, in this paper, all developed algorithms for tackling problem (P) will be based on the conditional gradient scheme with a unit step size called ConGradU. At this juncture, it is important to notice that Mangasarian [24] seems to have been the first work suggesting and analyzing the conditional gradient algorithm with a unit step size for maximizing a convex function over a polyhedron, in the context of machine learning problems.

A common and interesting appeal of the resulting algorithms is that they take the shape of a closed-form iterative scheme, i.e., they can be written as

$$x^{j+1} = \frac{\mathcal{S}(Ax^j)}{\|\mathcal{S}(Ax^j)\|_2}, \quad j = 0, 1, \dots$$

where  $\mathcal{S} : \mathbf{R}^n \rightarrow \mathbf{R}^n$  is a simple operator that can be either written in explicit form or efficiently computed<sup>3</sup>.

All problems addressed here are difficult nonconvex optimization problems and we make no claims with respect to global optimality; after all, these are difficult problems so obtaining cheap solutions must have some cost which here is a certificate of global optimality. Moreover, an important driving point is that there is no reason to discount stationary solutions of nonconvex problems versus globally optimal solutions to convex relaxations. For neither solution do we have a measure of the gap to the optimal solution of problem (1). We can empirically demonstrate similar solution quality, while the nonconvex methods are orders of magnitudes cheaper to compute and can be applied to data sets much larger than can be done with any known convex relaxation. This is in contrast to the well-known sparse recovery problem, for which there is an equivalence between the difficult combinatorial problem and the linear program relaxation when the data matrix satisfies certain conditions [8]. To put in perspective the development and results of this paper, we first discuss some of the relevant literature that has motivated this work.

## 1.1 Literature

The literature on sparse PCA can be divided according to the different modifications discussed above. Here, we briefly survey these approaches and detail them further in Section 5. With respect to the  $l_0$ -constrained PCA problem (1), thresholding [7] is perhaps the simplest approach, however is known to produce poor results. Sparse low-rank approximations (SLRA) of [39] looks for a more general  $(uv^T$  rather than  $xx^T$ ) sparse approximation by taking an approximation error level as input and determining the sparsity level  $k$  that is required to satisfy the desired error. Greedy methods [25] are also computationally expensive due to a maximum eigenvalue computation at each iteration, however an approximate greedy approach [10] offers a much cheaper way to derive an entire path of solutions (for each  $k = 0 \dots n$ ) which often suffice. These approaches are heuristics. The only globally optimal approach to this formulation is an exact search method [25] that is applicable only to extremely small problems.

The  $l_1$ -constrained PCA problem is a relaxation, or more precisely an upper bound, to problem (1). This problem was first considered in [18] and called SCoTLASS (simplified component technique least absolute shrinkage and selection). It was motivated by the LASSO (least absolute selection and shrinkage operator) approach used in statistics [33] for inducing sparsity in regression. In [34], a true penalty function is used to handle the  $l_1$  constraint and the resulting problem is solved as a system of differential equations (which requires a smooth approximation for the  $l_1$  constraint). While this approach was tested solely on a small 13-dimensional data set, we mention it as the only true penalty function approach in the sparse PCA literature.

---

<sup>3</sup>In fact, when  $\mathcal{S}$  is the identity operator, the scheme is nothing else but the well known power method to compute the first principal eigenvector of the matrix  $A$ , see, e.g., [32].

More recently, a computationally cheap approach to  $l_1$ -constrained PCA that can solve large-scale problems was given by [36]. We will show that this scheme is an application of the conditional gradient algorithm.

As already explained in the introduction, we are not interested in formulations that require expensive computations. Convex relaxations are either semidefinite-based or optimize over symmetric  $n \times n$  matrices, and are, indeed, much more computationally expensive than what we would like to consider here. Nevertheless, it is important to briefly recall some of these works. In [11], d’Aspremont et al. introduced a convex relaxation to  $l_1$ -constrained PCA using semidefinite programming, however this formulation can only be solved on very small dimensions ( $< 100$ ). This was the first approach with convex optimization to any sparse PCA modification and motivated a convex relaxation for  $l_1$ -penalized PCA for which better algorithms were given. Another approach for  $l_1$ -constrained PCA in [22] solves a different convex relaxation over  $\mathbf{S}^n$  based on a variational representation of the  $l_1$  ball (this relaxation turned out to be the dual of the semidefinite relaxation in [11]). This was the first convex relaxation for  $l_1$ -constrained PCA amenable to a medium number of dimensions (1000-2000).

We next turn to penalized sparse PCA where the sparsity-inducing term appears in the objective function. Several approaches are known for  $l_0$ -penalized PCA. The heuristic given in [39] is extended in [40] to this penalized version of PCA. In [10], the problem is reformulated as an equivalent convex maximization problem to which a semidefinite relaxation is applied, however, as mentioned above, solving such problems is too computationally expensive. More recently, [19] derived the same convex maximization problem (in a different manner) and proposed a first-order gradient-based algorithm which is identical to ConGradU. Another convex maximization representation was very recently presented in [31], whereby a specific parameterized concave approximation is used to replace the  $l_0$  term, and the resulting problem was solved by an iterative scheme called the minorization-maximization technique, which is, in fact, a specific instance of ConGradU.

The  $l_1$ -penalized PCA problem can be solved by a convex relaxation as shown in [11], but is amenable to only a medium number of dimensions. In [41], Zhou et al. considered a reformulation of PCA as a two variable regression problem to which an  $l_1$  penalty is added for one of the variables. While their approach is amenable to larger problems, it is not exactly  $l_1$ -penalized PCA and is still more computationally expensive (cf. Section 5.4) than other approaches we will discuss. Recently, [19] reformulated  $l_1$ -penalized PCA as an equivalent convex maximization problem (as with their  $l_0$ -penalized approach) that is also solved by the conditional gradient algorithm.

As discussed above, the conditional gradient algorithm has previously been applied to sparse PCA in various forms, so we now put the contributions of the above works into perspective. The works of [36, 29] detail their iterative schemes (without any general algorithm) specifically for sparse PCA. [31] details a general algorithm, similar to ConGradU, but meant for maximizing the difference of convex functions over a convex set. While maximizing a convex function is a subset of this algorithm, ConGradU, as detailed in Section 3, allows for nonconvex sets. [19] is the only previously known work that details a general algorithm for maximizing a convex function over a compact (possibly nonconvex) set. Indeed, the first-order algorithm proposed in [19] (labeled Algorithm 1 therein) is identical to what we term the ConGradU algorithm, however it is not recognized as the conditional gradient algorithm. As noticed in [19], both the  $l_0$  and  $l_1$ -penalized PCA algorithms they have proposed were earlier stated in [29], subject to slight modifications, who look for a general rank-one approximation (i.e.,  $uv^T$  rather than  $xx^T$ ); however, no convergence results were stated in [29].

## 1.2 Outline

We provide computationally simple approaches to both constrained and penalized versions of sparse PCA. It is important to recognize that all algorithms here are schemes for nonconvex problems; we pay the price of no global optimality criterion and gain in amenability to problem sizes that convex relaxations cannot handle.

In Section 2, we define the problems of interest and some of their properties. Section 3 recalls some basic optimality results for maximizing a convex function over a compact set. We then detail ConGradU, a specific conditional gradient scheme with unit step size, and establish its convergence properties. Section 4 provides a mathematical toolbox proving a series of propositions that are used to develop the known cheap algorithms mentioned above, as well as for deriving new schemes. These propositions are simple and easy to prove so we believe it benefits the reader to go through these tools first.

Section 5 then details the algorithms for all versions of sparse PCA. We start with a simple algorithm for the true  $l_0$ -constrained PCA problem (1). To the author's knowledge, this is the first available scheme that directly approaches this problem, is amenable to large-scale problems and proven to converge to a stationary point of problem (1). While the  $l_0$  constraint is a difficult nonsmooth and nonconvex constraint, we need not look for ways around this constraint, and rather we approach the given problem as is. The basis for our approach is the simple, yet surprising, result (cf. Section 4) that while maximizing a quadratic function over the  $l_2$  unit ball with an  $l_0$  constraint is a difficult problem, maximizing a linear function over the same nonconvex set is simple and can be solved in  $O(n)$  time. An important aspect of the new  $l_0$ -constrained PCA algorithm is that no parameters need be tuned in order to obtain a stationary point that has the exact desired sparsity<sup>4</sup>. The second main part of Section 5 focuses on all aforementioned iterative schemes which have been proposed in the literature. Building on the results of Section 4, we show that all these schemes can directly be obtained as a particular realization of ConGradU, or of some variant of it, thus providing a unifying framework to various seemingly different algorithmic approaches.

Section 6 provides experimental results and demonstrates the efficiency of many of the methods we have reviewed on large-scale problems. We show that they all give comparable solutions, i.e., very similar  $k$ -sparse solutions, with the advantage of  $l_0$ -constrained PCA being that the  $k$ -sparse solution is directly obtained at a lower computational cost. Section 7 ends with concluding remarks and briefly shows how to use the same tools to develop simple algorithms for related sparsity constrained problems.

## 1.3 Notation

We write  $\mathbf{S}^n$  ( $\mathbf{S}_+^n$ ,  $\mathbf{S}_{++}^n$ ) to denote the set of symmetric (positive-semidefinite, positive-definite) matrices of size  $n$  and  $\mathbf{R}^n$  ( $\mathbf{R}_+^n$ ,  $\mathbf{R}_{++}^n$ ) to denote the set of (nonnegative, strictly positive) real vectors of size  $n$ . The vector  $e$  is the  $n$ -vector of ones. Given a vector  $x \in \mathbf{R}^n$ ,  $\|x\|_2 = (\sum_i x_i^2)^{1/2}$  defines the  $l_2$  norm,  $\|x\|_0$  defines the cardinality of  $x$ , i.e., the number of nonzero entries of  $x$  and usually called here the  $l_0$  norm<sup>5</sup>, and  $\|x\|_\infty = \max(|x_1|, \dots, |x_n|)$ . Given a matrix  $X \in \mathbf{S}^n$ ,  $\|X\|_\infty = \max_{i,j} X_{i,j}$  and  $\|X\|_F = (\sum_{i,j} X_{i,j}^2)^{1/2}$ . For a vector  $x \in \mathbf{R}^n$ ,  $|x|$  denotes the vector with  $i^{\text{th}}$  entry  $|x_i|$ ,  $\text{sgn}(x)$  denotes the vector with  $i^{\text{th}}$  entry  $-1, 0, 1$  if  $x_i < 0, x_i = 0, x_i > 0$ , and  $x_+$  denotes the vector with  $i^{\text{th}}$  entry  $\max(x_i, 0)$ . For a vector  $x \in \mathbf{R}^n$ ,  $\text{diag}(x)$  denotes the diagonal matrix with  $x$  on its diagonal. For any nonzero integer  $n$ , denote the set  $\{1, \dots, n\}$  as  $[n]$ . Let  $I_n$  denote the identity matrix in dimension  $n$ . Let  $\mathcal{C}^1$  denote the space of

---

<sup>4</sup>All other algorithm based on modifications can be used to obtain a desired sparsity as well, however parameters must be tuned accordingly.

<sup>5</sup>We note an abuse of terminology because  $\|x\|_0$  is not a true norm since it is not positively homogenous.

once continuously differentiable functions on  $\mathbf{R}^n$ . Given an optimization problem (P), we use  $\text{argmax}(P)$  to denote its optimal solutions set.

## 2 Problem Formulations

This section describes the relationship between  $l_0$ -constrained PCA and the various modified sparse PCA problems that are discussed throughout the paper, as well as certain properties.

### 2.1 The Original Optimization Model

We start with some useful and elementary properties of the  $l_0$ -constrained PCA problem.

#### The $l_0$ -constrained PCA Problem

Given a symmetric matrix  $A \in S^n$  and sparsity level  $k \in [1, n]$ , the main problem of interest is to solve the  $l_0$ -constrained PCA problem (i.e., the sparse eigenvalue problem):

$$(E) \quad \max\{x^T A x : \|x\|_2 = 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\}. \quad (2)$$

In most applications,  $A$  is the covariance matrix of some data matrix  $B \in \mathbf{R}^{m \times n}$  such that  $A = B^T B$  and hence is positive semidefinite. The latter fact will be exploited in reformulations of the problem described below. In fact, as is very well-known, since problem (E) is constrained by the unit sphere, without loss of generality (see e.g., [10, 11, 19]) we can always assume that  $A$  is positive semidefinite since we clearly have

$$\max_{x \in \mathbf{R}^n} \{x^T A_\sigma x : \|x\|_2 = 1, \|x\|_0 \leq k\} = \max_{x \in \mathbf{R}^n} \{x^T A x : \|x\|_2 = 1, \|x\|_0 \leq k\} + \sigma,$$

with  $\sigma > 0$  such that  $A_\sigma := A + \sigma I_n \in S_{++}^n$ , i.e., with respect to optimal objective values we are solving the same problem.

The next result furthermore states that we can relax the sphere constraint to its convex counterpart, the unit ball, namely we consider the problem

$$(E_\sigma) \quad \max\{x^T A_\sigma x : \|x\|_2 \leq 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\}$$

and also show that problems (E) and  $(E_\sigma)$  are equivalent and admit the same set of optimal solutions. The simple proof is omitted.

**Lemma 1** *Fix any  $\sigma > 0$  such that  $A_\sigma \in S_{++}^n$ . Then,*

$$(a) \max\{x^T A_\sigma x : \|x\|_2 \leq 1, \|x\|_0 \leq k\} = \max\{x^T A x : \|x\|_2 = 1, \|x\|_0 \leq k\} + \sigma.$$

$$(b) \text{argmax}(E) = \text{argmax}(E_\sigma).$$

Needless to say that both problems (E) and  $(E_\sigma)$  remain hard nonconvex problems as they consist of maximizing a convex (strongly convex) function over a compact set, and clearly the possibility of using the convex relaxation  $\{x : \|x\|_2 \leq 1\}$  instead of the nonconvex unit sphere constraint does not change the situation. However, it is useful to know that either one of these constraints can be used when tackling the problem, and this will be done throughout the rest of the paper without any further mentioning.

Throughout the paper, (E) will be referred to as the  $l_0$ -constrained PCA problem, and  $A$  always denotes a symmetric matrix while  $A_\sigma$  will denote a symmetric positive definite matrix. We now consider the other formulations that will be analyzed.



## 2.2 Modified Optimization Models

As mentioned in the introduction, most approaches to sparse PCA solve one of the following modified problems. The first variation is a relaxation based on the relation

$$\|x\|_1 \leq \sqrt{\|x\|_0} \|x\|_2 \quad \forall x \in \mathbf{R}^n \quad (3)$$

which follows from the Cauchy-Schwarz inequality. The hard  $l_0$  constraint in problem (2) is replaced by an  $l_1$  constraint, resulting in

### The $l_1$ -constrained PCA Problem

$$\max \{x^T Ax : \|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}, x \in \mathbf{R}^n\}, \quad (4)$$

which thanks to inequality (3) is an upper bound to the original  $l_0$ -constrained PCA problem. Two other variations are based on penalizations of the  $l_0$  and  $l_1$  constraints of the above formulations. Note that the penalized terminology is different from the usual one used in optimization, and here is used to mention that it rather optimizes a tradeoff between how good and how sparse the approximation is. We first penalize the  $l_0$  constraint in problem (2), resulting in

### The $l_0$ -penalized PCA Problem

$$\max \{x^T Ax - s\|x\|_0 : \|x\|_2 \leq 1, x \in \mathbf{R}^n\}, \quad (5)$$

where  $s > 0$  is a parameter that must be tuned to achieve the truly desired sparsity level at which  $\|x\|_0 = k$ . However, to avoid the trivial optimal solution  $x^*(s) \equiv 0$ , the parameter  $s$  must be restricted. First recall the well-known norm relations

$$\|x\|_\infty \leq \|x\|_2 \leq \|x\|_1, \quad \forall x \in \mathbf{R}^n. \quad (6)$$

Using the Hölder inequality<sup>6</sup>, combined with inequalities (3) and (6), it follows that

$$x^T Ax - s\|x\|_0 \leq \|A\|_\infty \|x\|_1^2 - s\|x\|_0 \leq (\|A\|_\infty - s)\|x^*\|_0,$$

for all  $x$  feasible for problem (5). Thus, to avoid the trivial solution, it is assumed that  $s \in (0, \|A\|_\infty)$ . Note that while  $s \geq \|A\|_\infty$  necessarily implies the trivial solution, taking  $s < \|A\|_\infty$  does not guarantee we avoid it, but only gives a particular bound.

Likewise, a penalized version of the  $l_1$ -constrained PCA problem (4) yields

### The $l_1$ -penalized PCA Problem

$$\max \{x^T Ax - s\|x\|_1 : \|x\|_2 \leq 1, x \in \mathbf{R}^n\}. \quad (7)$$

Note that, as in problem (5), we need to restrict the value of the parameter  $s$  in order to avoid the trivial solution. Again, using the Hölder inequality and (6), it is easy to see that

$$x^T Ax - s\|x\|_1 \leq (\|A\|_F - s)\|x^*\|_1,$$

for all  $x$  feasible for problem (7), and hence it is assumed that  $s \in (0, \|A\|_F)$ .

---

<sup>6</sup>For any  $u, v \in \mathbf{R}^n$ , the Hölder inequality states that  $|\langle u, v \rangle| \leq \|u\|_p \|v\|_q$ , where  $p + q = pq, p \geq 1$ .

Again, note that the formulated penalized/relaxed problems remain hard nonconvex maximization problems despite the convexity of their constraints. In fact,  $l_0$ -penalized PCA and the  $l_1$ -penalized version share an additional difficulty in that their objectives are neither concave nor convex. In Sections 5.3 and 5.4, we show how this difficulty is overcome.

### The Approximate $l_0$ -penalized PCA Problem

The last approach for solving problem (2) involves approximating the  $l_0$  norm in the objective of  $l_0$ -penalized PCA. The idea of approximating the  $l_0$  norm by some nicer continuous functions naturally emerged from very well-known mathematical approximations of the step and sign functions (see, e.g., [6]). Indeed, it is easy to see that for any  $x \in \mathbf{R}^n$ , one can write

$$\|x\|_0 = \sum_{i=1}^n \text{sgn}(|x_i|).$$

Thus, formally, we want to replace the problematic expression  $\text{sgn}(|t|)$  by some nicer function and consider an approximation of the form

$$\|x\|_0 = \lim_{p \rightarrow 0} \sum_{i=1}^n \varphi_p(|x_i|),$$

where  $\varphi_p : \mathbf{R}_+ \rightarrow \mathbf{R}_+$  is an appropriately chosen smooth function. Here, we consider the class of smooth concave functions which are monotone increasing and normalized such that  $\varphi_p(0) = 0, \varphi_p'(0) > 0$ .

This suggests to approximate problem (5) by considering for  $p, s > 0$  the following approximate maximization problem:

$$\max \{x^T Ax - s \sum_{i=1}^n \varphi_p(|x_i|) : \|x\|_2 \leq 1, x \in \mathbf{R}^n\}. \quad (8)$$

Approximations of the  $l_0$  norm have been considered in various applied contexts, for instance in the machine learning literature, see, e.g., [24, 35].

There exist several possible choices for the function  $\varphi_p(\cdot)$  that approximate the step function, and for details the reader is referred to the classic book by Bracewell [6, Chapter 4]. For illustration, here we give the following examples as well as their graphical representation in Figure 1 (left).

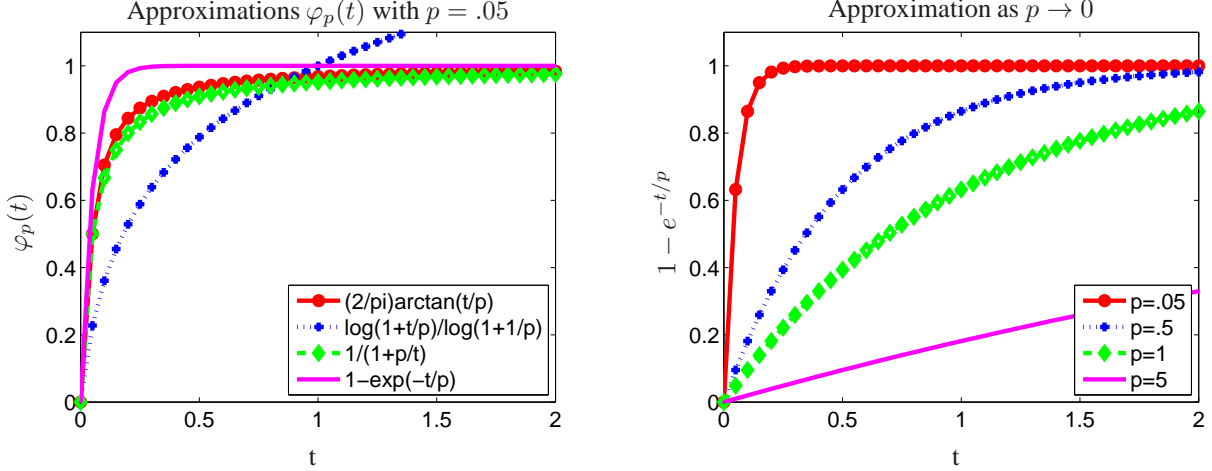
**Example 2** *Concave functions  $\varphi_p(\cdot), p > 0$*

- (a)  $\varphi_p(t) = (2/\pi) \tan^{-1}(t/p),$
- (b)  $\varphi_p(t) = \log(1 + t/p) / \log(1 + 1/p),$
- (c)  $\varphi_p(t) = (1 + p/t)^{-1},$
- (d)  $\varphi_p(t) = 1 - e^{-t/p}.$

The last example (d) was successfully used in the context of machine learning by Mangasarian [24], and gives

$$\varphi_p(|t|) = 1 - e^{-p|t|} = \begin{cases} 0 & \text{if } t = 0 \\ > 0 & \text{if } t \neq 0. \end{cases} \quad (9)$$





**Figure 1:** The left plot shows four concave functions  $\varphi_p(t)$  that can be used to approximate  $\|x\|_0 \approx \sum_{i=1}^n \varphi_p(|x_i|)$  for fixed  $p = .05$ . The right plot shows how the concave approximation  $1 - e^{-t/p}$  converges to the indicator function as  $p \rightarrow 0$ .

A nice feature of this example is that it also lower bounds the  $l_0$  norm, namely, we have

$$\sum_{i=1}^n \varphi_p(|x_i|) \leq \|x\|_0, \quad \forall x \in \mathbf{R}^n.$$

See Figure 1 (right) for its behavior for various values of  $p$ .

All problems listed in this section, as well as several other equivalent reformulations that will be derived in Section 5, will be solved by ConGradU, a conditional gradient algorithm with unit step size for maximizing a convex function over a compact set which is described next.

### 3 The Conditional Gradient Algorithm

#### 3.1 Background

The conditional gradient algorithm, is a well-known and simple gradient algorithm, see, e.g., [21, 12] and the book [4] for a general overview of this method as well as more references. Note that this algorithm dates back to 1956, and is also known as the Frank-Wolfe algorithm [15] that was originally proposed for solving linearly constrained quadratic programs.

The conditional gradient algorithm is presented here for maximization problems because of our interest in the sparse PCA problem. We first recall the conditional gradient algorithm for maximizing a continuously differentiable function  $F : \mathbf{R}^n \rightarrow \mathbf{R}$  over a nonempty compact convex set  $C \subset \mathbf{R}^n$ :

$$\max \{F(x) : x \in C\}. \quad (10)$$

The conditional gradient algorithm generates a sequence  $\{x^j\}$  via the iteration:

$$x^0 \in C, \quad x^{j+1} = x^j + \alpha^j(p^j - x^j), \quad j = 0, 1, \dots$$

where

$$p^j = \operatorname{argmax} \{ \langle x - x^j, \nabla F(x^j) \rangle : x \in C \}, \quad (11)$$

and where  $\alpha^j \in (0, 1]$  is a stepsize that can be determined by the Armijo or limited maximization rule [4]. It can be shown that every limit point of the sequence  $\{x^j\}$  generated by the conditional gradient algorithm is a stationary point. Furthermore, under various additional assumptions on the function  $F$  and/or the set  $C$  (e.g., strong convexity of the function  $F$  and/or of the set  $C$ ), rate of convergence results can be established, see in particular the work of Dunn [13] and references therein.

Clearly, the conditional gradient algorithm becomes attractive and simple when its main computational step (11) can be performed efficiently, e.g., when  $C$  is a bounded polyhedron it reduces to solving a linear program, or even better when it can be solved *analytically*.

### 3.2 Maximizing a Convex Function via ConGradU

As we shall see below, all potential reformulations of the sparse PCA problem will lead to maximizing a *convex* (possibly nonsmooth) function over a compact (possibly nonconvex) set  $C \subset \mathbf{R}^n$ . It will be shown that, for such a class of problems, the conditional gradient algorithm will reduce to a very simple iterative scheme. First we recall some basic definitions and properties relevant to maximizing convex functions.

Let  $F : \mathbf{R}^n \rightarrow \mathbf{R}$  be convex, and  $C \subset \mathbf{R}^n$  be a nonempty compact set. Throughout, we assume that  $F$  is not constant on  $C$ . We denote by  $F'(x)$  any subgradient of the convex function  $F$  at  $x$  which satisfies

$$F(v) - F(x) \geq \langle v - x, F'(x) \rangle, \quad \forall v \in \mathbf{R}^n. \quad (12)$$

The set of all subgradients of  $F$  at  $x$  is the subdifferential of the function  $F$  at  $x$ , denoted by  $\partial F(x)$ , i.e.,

$$\partial F(x) = \{g : F(v) \geq F(x) + \langle v - x, g \rangle, \quad \forall v \in \mathbf{R}^n\},$$

which is a closed convex set. When  $F \in \mathcal{C}^1$ , the subdifferential reduces to a singleton which is the gradient of  $F$ , that is  $\partial F(x) = \{\nabla F(x)\}$ , and (12) is the usual gradient inequality for the convex function  $F$ . In the following, we use the notation  $F'(\cdot)$  to refer to either a gradient or subgradient of  $F$ ; the context will be clear in the relevant situation.

The first result recalls two useful properties when maximizing a convex function (see [28, Section 32]).

**Proposition 3** *Let  $F : \mathbf{R}^n \rightarrow \mathbf{R}$  be convex,  $S \subset \mathbf{R}^n$  be an arbitrary set and let  $\operatorname{conv}(S)$  denote its convex hull. Then,*

(a)  $\sup\{F(x) : x \in \operatorname{conv}(S)\} = \sup\{F(x) : x \in S\}$ , where the first supremum is attained only if the second is attained.

(b) If  $S \subset \mathbf{R}^n$  is closed with nonempty boundary  $\operatorname{bd}(S)$ , and  $F$  is bounded above on  $S$ , then  $\sup\{F(x) : x \in S\} = \sup\{F(x) : x \in \operatorname{bd}(S)\}$ .

The next result gives a first-order optimality criterion for maximizing a convex function  $F$  over a compact set  $C \subset \mathbf{R}^n$ . It uses property (a) of Proposition 3 and follows from [28, Corollary 32.4.1].

**Proposition 4** *Let  $F : \mathbf{R}^n \rightarrow \mathbf{R}$  be convex. If  $x$  is a local maximum of  $F$  over the nonempty compact set  $C$ , then*

$$(FOC) \quad \langle v - x, F'(x) \rangle \leq 0, \quad \forall v \in C. \quad (13)$$

When  $F \in \mathcal{C}^1$ , a point  $x \in C$  satisfying the first-order optimality criteria (FOC) (13) will be referred as a stationary point, otherwise, when  $F$  is convex nonsmooth, we will say that the point  $x \in C$  satisfies (FOC).

We are now ready to state the algorithm. It turns out that when the function  $F$  is convex and quadratic, we can also eliminate the need for finding a step size and consider the conditional gradient algorithm with a fixed unit step size and still preserve its convergence properties. This simplified version of the conditional gradient algorithm will be used throughout the paper and referred to as ConGradU.

---

**Algorithm 1 ConGradU – Conditional Gradient Algorithm with Unit Step Size**

---

**Require:**  $x^0 \in C$

- 1:  $j \leftarrow 0$ ,
  - 2: **while** stopping criteria **do**
  - 3:    $x^{j+1} \in \operatorname{argmax}\{\langle x - x^j, F'(x^j) \rangle : x \in C\}$
  - 4: **end while**
  - 5: **return**  $x^j$
- 

To analyze ConGradU, in what follows, it will be convenient to introduce the quantity

$$\gamma(x) := \max\{\langle v - x, F'(x) \rangle : v \in C\} \quad (14)$$

for any  $x \in \mathbf{R}^n$ . Since  $C$  is compact, this quantity is well-defined and admits a global maximizer

$$u(x) \in \operatorname{argmax}\{\langle v - x, F'(x) \rangle : v \in C\},$$

and thus, we have  $\gamma(x) = \langle u(x) - x, F'(x) \rangle$ .

In terms of the above defined quantities, we thus have that  $x^*$  satisfies (FOC) is equivalent to saying that  $x^*$  is a global maximizer of  $\psi(v) = \langle v - x^*, F'(x^*) \rangle$ , i.e.,

$$x^* \in \operatorname{argmax}\{\langle v - x^*, F'(x^*) \rangle : v \in C\} = u(x^*).$$

Thus the ConGradU algorithm is nothing else but a fixed point scheme for the map  $u(\cdot)$ , and simply reads as

$$x^0 \in C, \quad x^{j+1} = u(x^j), \quad j = 0, 1, \dots$$

**Lemma 5** *Let  $F : \mathbf{R}^n \rightarrow \mathbf{R}$  be convex,  $C \subset \mathbf{R}^n$  be nonempty and compact, and let  $\gamma(x)$  be given by (14). Then,*

- (i)  $\gamma(x) \geq 0$  for all  $x \in C$ .
- (ii) For any  $v \in C$  and any  $x \in \mathbf{R}^n$ ,

$$F(u(x)) - F(x) \geq \gamma(x) \geq \langle v - x, F'(x) \rangle.$$

**Proof.** The proof of (i) and the right inequality of (ii) follows immediately from the definition of  $\gamma(x)$ , while the left inequality in (ii) follows from the subgradient inequality for the convex function  $F$ :

$$F(u(x)) - F(x) \geq \langle u(x) - x, F'(x) \rangle = \gamma(x).$$

■

We are ready to state the convergence properties of ConGradU.

**Theorem 6** Let  $F : \mathbf{R}^n \rightarrow \mathbf{R}$  be a convex function and let  $C \subset \mathbf{R}^n$  be nonempty compact. Let  $\{x^j\}$  be the sequence generated by the algorithm ConGradU. Then the following statements hold:

(a) the sequence of function values  $F(x^j)$  is monotonically increasing and

$$\lim_{j \rightarrow \infty} \gamma(x^j) = 0.$$

(b) If for some  $j$  the iterate  $x^j$  satisfies  $\gamma(x^j) = 0$ , then the algorithm stops with  $x^j$  satisfying (FOC). Otherwise the algorithm generates an infinite sequence  $\{x^j\}$  with strictly increasing function values  $\{F(x^j)\}$ .

(c) Moreover, if  $F$  is continuously differentiable, then every limit point of the sequence  $\{x^j\}$  converges to a stationary point.

**Proof.** By definition of the iteration of ConGradU, using our notations, the sequence  $\{x^j\}$  is well defined via  $x^{j+1} \in u(x^j), \forall j = 0, 1, \dots$ . Invoking Lemma 5, we obtain

$$0 \leq \gamma(x^j) \leq F(x^{j+1}) - F(x^j), \quad \forall j = 0, 1, \dots,$$

showing that the sequence  $\{F(x^j)\}$  is monotone increasing. Summing the inequality above for  $j = 0, \dots, N-1$ , we have

$$\sum_{j=0}^{N-1} \gamma(x^j) \leq F(x^N) - F(x^0).$$

Since  $C$  is compact and  $F(\cdot)$  is continuous, we also have  $F(x^N) \leq \max\{F(x) : x \in C\} := F_*$ , and thus it follows from the above inequality that  $\sum_{j=0}^{N-1} \gamma(x^j) \leq F_* - F(x^0)$ , and hence the nonnegative series  $\sum_{j=0}^{\infty} \gamma(x^j)$  is convergent so that  $\gamma(x^j)$  converges to 0. Now, since  $\gamma(x^j) \geq 0$  for all  $j = 0, \dots$ , then if for some  $j$  the iterate  $x^j$  is such that  $\gamma(x^j) = 0$  then the procedure stops at iteration  $j$  with  $x^j$  satisfying (FOC). Otherwise,  $\gamma(x^j) > 0$  and the iteration generates an infinite sequence with  $F(x^{j+1}) > F(x^j)$ . In the latter case, assuming now that  $F \in C^1$ , we will now prove the last statement of the theorem. Since  $C$  is compact, the sequence  $\{x^j\} \subset C$  is bounded. Passing to subsequences if necessary, for any limit point  $x^\infty$  of  $\{x^j\}$  we thus have  $x^j \rightarrow x^\infty$ . Without loss of generality we let  $u(x^j) \rightarrow \bar{u}$ . Using the facts

$$\gamma(x^j) = \langle u(x^j) - x^j, F'(x^j) \rangle \quad \text{and} \quad \langle v - x^j, F'(x^j) \rangle \leq \langle u(x^j) - x^j, F'(x^j) \rangle, \quad \forall v \in C,$$

and since  $\gamma(x^j) \rightarrow 0$  and  $F \in C^1$ , passing to the limit over appropriate subsequences in the above relations, it follows that  $\langle \bar{u} - x^\infty, F'(x^\infty) \rangle = 0$  and  $\langle v - x^\infty, F'(x^\infty) \rangle \leq \langle \bar{u} - x^\infty, F'(x^\infty) \rangle, \forall v \in C$ , and hence  $\langle v - x^\infty, F'(x^\infty) \rangle \leq 0, \forall v \in C$ , showing that  $x^\infty$  is stationary. ■

**Remark 7** (a) For the case of convex  $F$  and bounded polyhedron  $C$ , as noted in the introduction, Mangasarian [24] seems to have been the first work suggesting the possibility of using a unitary step size in the conditional gradient scheme and proved that the algorithm generates a finite sequence (thanks to the polyhedrality of  $C$ ) that terminates at a stationary point.

(b) Very recently, the use of a unitary stepsize in the conditional gradient scheme was rediscovered in [19] with  $C$  being an arbitrary compact set, which is identical to Algorithm 1.

(c) The proof of Theorem 6 is patterned after the one given in [24]. Note that part (a) also follows from [19] as well. Furthermore, under stronger assumptions on  $F$  and  $C$ , [19, Theorem 4] also established a stepsize convergence rate giving an upper estimate on the number of iterations the algorithm takes to produce a step of small size.

(d) Finally, as kindly pointed out by a referee, the proof of convergence could also probably be derived by using the general approach developed in the classical monograph of Zangwill [37].

We end this section with a particular realization of ConGradU for an interesting class of problems given as

$$(G) \quad \max_x \{f(x) + g(|x|) : x \in C\}$$

where

$$\begin{aligned} f : \mathbf{R}^n &\rightarrow \mathbf{R} && \text{is convex,} \\ g : \mathbf{R}_+^n &\rightarrow \mathbf{R} && \text{is convex differentiable and monotone decreasing}^7, \\ C \subseteq \mathbf{R}^n &&& \text{is a compact set.} \end{aligned}$$

Our interest for this form is mainly driven by and is particularly useful for handling the case of the approximate  $l_0$ -penalized problem (cf. Section 5.3), which precisely has the form of this optimization model with adequate choice of the kernel  $\varphi_p$  that can be used to approximate the  $l_0$  norm (cf. Section 2). It will also allow for developing a novel simple scheme for the  $l_1$ -penalized problem (cf. Section 5.4).

Note that under the stated assumptions for  $g$ , the composition  $g(|x|)$  is not necessarily convex and thus ConGradU cannot be applied to problem (G). However, thanks to the componentwise monotonicity of  $g(\cdot)$ , it is easy to see that problem (G) can be recast as an equivalent problem with additional constraints and variables as follows:

$$(GG) \quad \max_{x,y} \{f(x) + g(y) : |x| \leq y, x \in C\}.$$

Note that, without loss of generality, an additional upper bound can be imposed on  $y$  in order to enforce compactness of the feasible region of problem (GG) (e.g., by setting the upper bound on  $y$  as  $y_c := \arg\max\{|x| : x \in C\}$ ), but this need not be computed in order to establish our following result. Clearly, the new objective of (GG) is now convex in  $(x, y)$  and thus we can apply ConGradU. We show below that the main iteration in that case leads to an attractive weighted  $l_1$ -norm maximization problem, which, in turn, as shown in Section 4 can be solved in closed form for the cases of interest, namely when  $C$  is the compact set described by the unit sphere or unit ball.

**Proposition 8** *The algorithm ConGradU applied to problem (GG) generates a sequence  $\{x^j\}$  by solving the weighted  $l_1$ -norm maximization problem*

$$x^0 \in C, \quad x^{j+1} = \arg\max \left\{ \langle a^j, x \rangle - \sum_i w_i^j |x_i| : x \in C \right\}, j = 0, \dots, \quad (15)$$

where  $w^j = -g'(|x^j|) > 0$  and  $a^j = f'(x^j) \in \mathbf{R}^n$ .

**Proof.** Applying ConGradU to the convex function  $H(x, y) := f(x) + g(y)$ , with  $y^0 = |x^0|, x^0 \in C$  we obtain,

$$\begin{aligned} (x^{j+1}, y^{j+1}) &= \arg\max_{x,y} \{ \langle x - x^j, f'(x^j) \rangle + \langle y - y^j, g'(y^j) \rangle : x \in C, |x| \leq y \} \\ &= \arg\max_{x \in C} \left\{ \langle x - x^j, f'(x^j) \rangle + \max_y \{ \langle y - y^j, g'(y^j) \rangle : |x| \leq y \} \right\} \\ &= \arg\max_{x \in C} \{ \langle x - x^j, f'(x^j) \rangle + \langle |x| - y^j, g'(y^j) \rangle \}, \end{aligned}$$

---

<sup>7</sup>By monotone decreasing, we mean componentwise, i.e., each component of  $g'(\cdot)$  the gradient of  $g$  is such that  $g'_i(\cdot)_i < 0 \quad \forall i \in [n]$ .

where the last max computation with respect to  $y$  uses the fact that  $g'$  is monotone decreasing. It is also clear that, given the initialization  $y^0 = |x^0|$  and  $y^{j+1} = \operatorname{argmax} \{ \langle y, g'(y^j) \rangle : |x^{j+1}| \leq y \}$ , we have  $y^j = |x^j|$  for all  $j = 0, 1, \dots$ . Omitting constant terms, the last iteration can be simply rewritten as

$$x^{j+1} = \operatorname{argmax}_{x \in C} \{ \langle x, f'(x^j) \rangle + \langle |x|, g'(|x^j|) \rangle \},$$

which with  $w^j := -g'(|x^j|) > 0$  proves the desired result stated in (15). ■

## 4 A Simple Toolbox for Building Simple Algorithms

The algorithms discussed in this paper are based on the fact that the  $l_0$ -constrained PCA problem as well as many of the penalized and constrained  $l_0$  and  $l_1$  optimization problems that are solved by the proposed conditional gradient algorithm have iterations that either have closed form solutions or are easily solvable. Specifically, the main step of ConGradU maximizes a linear function over a compact set, and the following lemmas and propositions show that for certain compact sets (e.g.  $\{x \in \mathbf{R}^n : \|x\|_2 = 1, \|x\|_0 \leq k\}$  and  $\{x \in \mathbf{R}^n : \|x\|_2 = 1, \|x\|_1 \leq k\}$ ), these subproblems are easy to solve.

In addition, propositions are given that will be used to reformulate (in Section 5) problems such as  $l_0$  and  $l_1$ -penalized PCA, which have neither convex nor concave objectives, to problems that maximize a convex objective function. The propositions for maximizing linear functions over compact sets can then be used in ConGradU as applied to the reformulated  $l_0$  and  $l_1$ -penalized PCA problems. Combined with the conditional gradient algorithm discussed above, these propositions are the only tools required throughout the remainder of the paper.

We start with an obvious but very useful lemma that is used in all of the following propositions.

**Lemma 9** *Given  $0 \neq a \in \mathbf{R}^n$ ,*

$$\max \{ \langle a, x \rangle : \|x\|_2 = 1, x \in \mathbf{R}^n \} = \max \{ \langle a, x \rangle : \|x\|_2 \leq 1, x \in \mathbf{R}^n \} = \|a\|_2,$$

*with maximizer  $x^* = a/\|a\|_2$ .*

**Proof.** Immediate from Cauchy-Schwarz inequality. ■

The following propositions make use of the following operator:

**Definition 10** *Given any  $a \in \mathbf{R}^n$ , define the operator  $T_k : \mathbf{R}^n \rightarrow \mathbf{R}^n$  by*

$$T_k(a) := \operatorname{argmin}_x \{ \|x - a\|_2^2 : \|x\|_0 \leq k, x \in \mathbf{R}^n \}.$$

This operator is thus the best  $k$ -sparse approximation of a given vector  $a$ . Despite the nonconvexity of the constraint, it is easy to see that  $(T_k(a))_i = a_i$  for the  $k$  largest entries (in absolute value) of  $a$  and  $(T_k(a))_i = 0$  otherwise. In case the  $k$  largest entries are not uniquely defined, we select the smallest possible indices.

In other words, without loss of generality, with  $a \in \mathbf{R}^n$  such  $|a_1| \geq \dots \geq |a_n|$ , we have

$$(T_k(a))_i = \begin{cases} a_i, & i \leq k; \\ 0, & \text{otherwise.} \end{cases}$$



Computing  $T_k(\cdot)$  only requires determining the  $k^{th}$  largest number of a vector of  $n$  numbers which can be done in  $O(n)$  time [5] and zeroing out the proper components in one more pass of the  $n$  numbers.

The next proposition is an extension of Lemma 9 to  $l_0$ -constrained problems. This is the simple result that maximizing a linear function over the nonconvex set  $\{x \in \mathbf{R}^n : \|x\|_2 = 1, \|x\|_0 \leq k\}$  is equally simple and can be solved in  $O(n)$  time.

**Proposition 11** *Given  $0 \neq a \in \mathbf{R}^n$ ,*

$$\max_x \{\langle a, x \rangle : \|x\|_2 = 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\} = \|T_k(a)\|_2 \quad (16)$$

*with solution obtained at*

$$x^* = \frac{T_k(a)}{\|T_k(a)\|_2}.$$

**Proof.** By Lemma 9, the optimal value of problem (16) is  $\sqrt{\sum_{i \in \mathcal{I}} a_i^2}$  for some subset of indices  $\mathcal{I} \subseteq [n]$  with  $|\mathcal{I}| \leq k$ . The set  $\mathcal{I}$  that maximizes this value clearly contains the indices of the  $k$  largest elements of the vector  $|a|$ . Thus, by definition of  $T_k(a)$ , solving problem (16) is equivalent to solving

$$\max_x \{\langle x, T_k(a) \rangle : \|x\|_2 = 1, x \in \mathbf{R}^n\}$$

from which the result follows by Lemma 9. ■

Another version of this result gives the solution with a squared objective.

**Proposition 12** *Given  $a \in \mathbf{R}^n$ ,*

$$\max_x \{\langle a, x \rangle^2 : \|x\|_2 = 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\} = \|T_k(a)\|_2^2$$

*with solution obtained at*

$$x^* = \frac{T_k(a)}{\|T_k(a)\|_2}.$$

**Proof.** Notice that the optimal objective value in Proposition 11 is nonnegative, and hence the squared objective value here does not change the optimal solution to the problem with linear objective. The result follows. ■

The above propositions will serve to derive the novel schemes given in Section 5. The next result is a modification of Proposition 12. It shows that the  $l_0$ -penalized version of maximizing a squared linear function yields a closed-form solution.

**Proposition 13** *Given  $a \in \mathbf{R}^n, s > 0$ ,*

$$\max_x \{\langle a, x \rangle^2 - s\|x\|_0 : \|x\|_2 = 1, x \in \mathbf{R}^n\} = \sum_{i=1}^n (a_i^2 - s)_+$$

*is solved by*

$$x_i^* = \frac{a_i [\text{sgn}(a_i^2 - s)]_+}{\sqrt{\sum_{j=1}^n a_j^2 [\text{sgn}(a_j^2 - s)]_+}}.$$

**Proof.** Assume without loss of generality that  $|a_1| \geq \dots \geq |a_n|$ . The problem can be rewritten as

$$\max_{p \in [n]} \{-sp + \max_x \{\langle a, x \rangle^2 : \|x\|_2 = 1, \|x\|_0 \leq p, x \in \mathbf{R}^n\}\}.$$

Using Proposition 12, the inner maximization in  $x$  is solved at

$$x_i^* = \begin{cases} a_i / \sqrt{\sum_{j=1}^p a_j^2}, & i \leq p; \\ 0, & \text{otherwise,} \end{cases}$$

and the problem is equal to

$$\max_{p \in [n]} \{-sp + \|T_p(a)\|_2^2\} = \max_{p \in [n]} \left\{ \sum_{i=1}^p a_i^2 - sp \right\} = \max_{p \in [n]} \left\{ \sum_{i=1}^p (a_i^2 - s) \right\} = \sum_{i=1}^n (a_i^2 - s)_+.$$

Notice that the optimal  $p$  is the largest index  $i$  such that  $a_i^2 \geq s$  which makes the above expression for  $x^*$  equivalent to the expression in the proposition. ■

Our next two results are concerned with  $l_1$ -penalized/constrained optimization problems. First, it is useful to recall the following well-known operators which are particular instances of the so-called Moreau's proximal map [26]; see, for instance, [9] for these results and many more. Given  $a \in \mathbf{R}^n$  and  $W$  an  $n \times n$  diagonal matrix  $W = \text{diag}(w)$ ,  $w \in \mathbf{R}^n$  with positive entries  $w_i$ , let

$$\|Wx\|_1 := \sum_{i=1}^n w_i |x_i|; \quad \mathbb{B}_\infty^w := \{x \in \mathbf{R}^n : \|W^{-1}x\|_\infty \leq 1\}.$$

Then,

$$S_w(a) := \underset{x}{\operatorname{argmin}} \left\{ \frac{1}{2} \|x - a\|_2^2 + \|Wx\|_1 \right\} = (|a| - w)_+ \operatorname{sgn}(a), \quad (17)$$

$$\Pi_{\mathbb{B}_\infty^w}(a) := \underset{x}{\operatorname{argmin}} \{ \|x - a\|_2 : x \in \mathbb{B}_\infty^w \} = \operatorname{sgn}(a) \min\{w, |a|\} = a - S_w(a), \quad (18)$$

where  $S_w(a)$  and  $\Pi_{\mathbb{B}_\infty^w}(a)$  are respectively known as the soft-thresholding operator and the projection operator.

**Proposition 14** For  $a \in \mathbf{R}^n$ ,  $w \in \mathbf{R}_{++}^n$ , and  $W = \text{diag}(w)$

$$\max \{ \langle a, x \rangle - \|Wx\|_1 : \|x\|_2 \leq 1, x \in \mathbf{R}^n \} = \sqrt{\sum_{i=1}^n (|a_i| - w_i)_+^2} = \|S_w(a)\|$$

which is solved by

$$x^* = S_w(a) / \|S_w(a)\|_2.$$

**Proof.** By the Hölder inequality, we have  $\|Wx\|_1 = \max_{\|v\|_\infty \leq 1} \langle v, Wx \rangle = \max\{ \langle z, x \rangle : z \in \mathbb{B}_\infty^w \}$ . Using the latter, we obtain

$$\begin{aligned} \max\{ \langle a, x \rangle - \|Wx\|_1 : \|x\|_2 \leq 1 \} &= \max_{\|x\|_2 \leq 1} \min_{z \in \mathbb{B}_\infty^w} \langle a - z, x \rangle \\ &= \min_{z \in \mathbb{B}_\infty^w} \max_{\|x\|_2 \leq 1} \langle a - z, x \rangle \\ &= \min\{ \|a - z\|_2 : z \in \mathbb{B}_\infty^w \} = \|S_w(a)\|_2, \end{aligned}$$

where the second equality follows from standard min-max duality [28], the third from Lemma 9, and the last one from using the relations (17)-(18), where the optimal  $z^* = a - S_w(a)$  and  $x^*$  follows from Lemma 9. ■

We next turn to the  $l_2/l_1$ -constrained problem,

$$\max\{\langle a, x \rangle : \|x\|_2 \leq 1, \|x\|_1 \leq k, x \in \mathbf{R}^n\}, \quad (19)$$

and state a result similar to Proposition 11 for maximizing a linear function over the intersection of the  $l_2$  unit ball with an  $l_1$  constraint. While maximizing over the intersection of the  $l_2$  unit ball with an  $l_0$  ball has an analytic solution, here we need an additional simple one dimensional search to express the solution of (19) via its dual.

**Proposition 15** *Given  $a \in \mathbf{R}^n$ , we have*

$$\max\{\langle a, x \rangle : \|x\|_2 \leq 1, \|x\|_1 \leq k, x \in \mathbf{R}^n\} = \min_{\lambda \geq 0} \{\lambda k + \|S_{\lambda e}(a)\|_2\}, \quad (20)$$

*the right hand side being a dual of (19). Moreover, if  $\lambda^*$  solves the one-dimensional dual, then an optimal solution of (19) is given by  $x^*(\lambda^*)$  where:*

$$x^*(\lambda) = S_{\lambda e}(a) / \|S_{\lambda e}(a)\|_2, \quad (e \equiv (1, \dots, 1) \in \mathbf{R}^n). \quad (21)$$

**Proof.** Dualizing only the  $l_1$  constraint, standard Lagrangian duality [28] implies:

$$\max\{\langle a, x \rangle : \|x\|_2 \leq 1, \|x\|_1 \leq k, x \in \mathbf{R}^n\} = \min\{\lambda k + \psi(\lambda) : \lambda \geq 0\},$$

with

$$\psi(\lambda) := \max_{\|x\|_2 \leq 1} \{\langle a, x \rangle - \lambda \|x\|_1\} = \|S_{\lambda e}(a)\|_2$$

where the last equality follows from Proposition 14 with  $x^*(\lambda)$  as given in (21). ■

The above propositions will be used to create simple algorithms for  $l_0$  and  $l_1$ -constrained and penalized PCA.

## 5 Sparse PCA via Conditional Gradient Algorithms

This section details algorithms for solving the original  $l_0$ -constrained PCA problem (2) and its three modifications (4), (5), and (7). Everything is developed using the simple tool box from Section 4. We derive novel algorithms as well as other known schemes that are shown to be particular realizations of the conditional gradient algorithm. A common and interesting appeal of all these algorithms is that they take the shape of a generalized power method, i.e., they can be written as

$$x^{j+1} = \frac{\mathcal{S}(Ax^j)}{\|\mathcal{S}(Ax^j)\|_2}$$

where  $\mathcal{S} : \mathbf{R}^n \rightarrow \mathbf{R}^n$  is a simple operator given in closed form or that can be computed very efficiently.

Table 1 summarizes the different algorithms that are discussed throughout this section. Except for the alternating minimization algorithm for  $l_1$ -penalized PCA of [41], they are all particular realizations of the ConGradU algorithm with an appropriate choice of the objective/constraints  $(F, C)$ .

Type	Iteration	Per-Iteration Complexity	References
$l_0$ -constrained	$x^j = \frac{T_k((A + \frac{\sigma}{2}I_n)x^j)}{\ T_k((A + \frac{\sigma}{2}I_n)x^j)\ _2}$	$O(kn), O(mn)$	novel
$l_1$ -constrained	$x_i^{j+1} = \frac{\text{sgn}(((A + \frac{\sigma}{2})x^j)_i)( ((A + \frac{\sigma}{2})x^j)_i  - \lambda^j)_+}{\sqrt{\sum_h ( ((A + \frac{\sigma}{2})x^j)_h  - \lambda^j)_+^2}}$	$O(n^2), O(mn)$	[36]
$l_1$ -constrained	$x_i^{j+1} = \frac{\text{sgn}((Ax^j)_i)( (Ax^j)_i  - s^j)_+}{\sqrt{\sum_h ( (Ax^j)_h  - s^j)_+^2}}$ where $s^j$ is $(k + 1)$ -largest entry of vector $ Ax^j $	$O(n^2), O(mn)$	[30]
$l_0$ -penalized	$z^{j+1} = \frac{\sum_i [\text{sgn}((b_i^T z^j)^2 - s)_+ (b_i^T z^j) b_i]}{\ \sum_i [\text{sgn}((b_i^T z^j)^2 - s)_+ (b_i^T z^j) b_i]\ _2}$	$O(mn)$	[29, 19]
$l_0$ -penalized	$x_i^{j+1} = \frac{\text{sgn}(2(Ax^j)_i)( 2(Ax^j)_i  - s\varphi'_p( x_i^j ))_+}{\sqrt{\sum_h ( 2(Ax^j)_h  - s\varphi'_p( x_h^j ))_+^2}}$	$O(n^2)$	[31]
$l_1$ -penalized	$y^{j+1} = \text{argmin}_y \{\sum_i \ b_i - x^j y^T b_i\ _2^2 + \lambda \ y\ _2^2 + s \ y\ _1\}$ $x^{j+1} = \frac{(\sum_i b_i b_i^T) y^{j+1}}{\ (\sum_i b_i b_i^T) y^{j+1}\ _2}$	See Section 5.4	[41]
$l_1$ -penalized	$x_i^{j+1} = \frac{\text{sgn}(((A + \frac{\sigma}{2})x^j)_i)( ((A + \frac{\sigma}{2})x^j)_i  - s)_+}{\sqrt{\sum_h ( ((A + \frac{\sigma}{2})x^j)_h  - s)_+^2}}$	$O(n^2), O(mn)$	novel
$l_1$ -penalized	$z^{j+1} = \frac{\sum_i ( b_i^T z^j  - s)_+ \text{sgn}(b_i^T z^j) b_i}{\ \sum_i ( b_i^T z^j  - s)_+ \text{sgn}(b_i^T z^j) b_i\ _2}$	$O(mn)$	[29, 19]

**Table 1:** Sparse PCA Algorithms. For each iteration,  $B \in \mathbf{R}^{m \times n}$  is a data matrix with  $A = B^T B$ .  $b_i$  is the  $i^{th}$  column of  $B$  (except for the  $l_1$ -penalized PCA of [41] where it is the  $i^{th}$  row of  $B$ ). For iterations with two complexities, the first uses the covariance matrix  $A$  and the second uses the decomposition  $A = B^T B$  to compute matrix-vector products as  $Ax$  or  $B^T(Bx)$ . Several iterations have two complexities, depending on whether data matrix  $B$  is available. The regularized  $l_1$ -constrained version of [36] is also novel. The  $l_0$  and  $l_1$ -penalized iterations of [29] require an  $O(mn)$  transformation to recover a sparse  $x^*$  from  $z^*$ .

We stress that the first algorithm for  $l_0$ -constrained PCA is the only algorithm that applies directly to the original and unmodified  $l_0$ -constrained PCA problem. The other algorithms are applied to modified problems where tuning a parameter is needed to get an approximation to the desired  $k$ -sparse problem. Unless otherwise specified,  $A$  is only assumed to be symmetric and  $A_\sigma$  is used to denote the regularized positive definite matrix. The exceptions are only when we assume we are given a data matrix  $B$  so that

$$A = B^T B.$$

## 5.1 $l_0$ -Constrained PCA

In this section, we focus on the original  $l_0$ -constrained PCA problem

$$\max \{x^T A x : \|x\|_2 = 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\}. \quad (22)$$

We apply the ConGradU algorithm with the constraint set  $C = \{x \in \mathbf{R}^n : \|x\|_2 = 1, \|x\|_0 \leq k\}$  and the convex objective (cf. Section 2):

$$q_\sigma(x) = x^T (A + \sigma I) x = x^T A_\sigma x, \sigma \geq 0.$$

When  $A \in S^n$  is already given positive semidefinite, the objective is already convex and there is no need for regularization and thus we simply set  $\sigma = 0$ . The resulting main iteration of ConGradU reduces to,

$$x^{j+1} = \operatorname{argmax} \{\langle x, A_\sigma x^j \rangle : \|x\|_2 = 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\} = \frac{T_k(A_\sigma x^j)}{\|T_k(A_\sigma x^j)\|_2}, j = 0, 1, \dots \quad (23)$$

with  $x^0 \in C$  and where the second equality is due to Proposition 11.

This novel iteration is obtained by maximizing a continuously differentiable convex function over a compact nonconvex set, and by Theorem 6, every one of its limit points converges to a stationary point. The complexity of each iteration requires computing  $A_\sigma x^j$  where  $A_\sigma \in \mathbf{S}^{n \times n}$  and  $x^j$  is  $k$ -sparse so the matrix-vector product requires  $O(nk)$  complexity. Computing the  $T_k(\cdot)$  operator is  $O(n)$  so each iteration is  $O(nk)$ . For very large problems where only a data matrix  $B \in \mathbf{R}^{m \times n}$  can be stored, the matrix-vector product  $A_\sigma x^j$  is computed as  $B^T(Bx^j)$ , which requires complexity  $O(mn)$ , so each iteration is actually  $O(mn)$ . The new scheme based on iteration (23) is an extremely simple way to approach the original  $l_0$ -constrained PCA problem (22) for any given matrix  $A_\sigma \in S^n$  and is the cheapest known non-greedy approach to  $l_0$ -constrained PCA.

Thus, a very simple gradient algorithm (and not greedy heuristics) can be directly applied to the desired problem. To put our novel scheme in perspective, let us recall the work [25] which offers the following greedy scheme. Given a set of  $k$  indices for variables with nonzero values,  $n - k$  subsets of indices are created by appending the  $k$  indices with one from the  $n - k$  remaining indices. Then  $n - k$  possible matrices are computed from the  $n - k$  groups of indices and the matrix with the maximum eigenvalue gives the index that provides the  $k + 1$ -sparse PCA solution. A full path of solutions can be computed for all values of  $k$ , but at a costly expense  $O(n^4)$  (or up to  $k$ -sparsity in  $O(kn^3)$ ), and with no statements on stationarity that we have. The work [25] also produces a branch-and-bound method for computing the exact solution, however, this method is amenable to only very small problems. The authors of [10] considered what they call an *approximate* greedy algorithm that generates an entire path of solutions up to  $k$ -sparsity in  $O(kn^2)$  ( $k \in \{1, \dots, n\}$ ). Rather than computing the  $n - k$  maximum eigenvalues each iteration, their scheme computes  $n - k$  dot products each iteration. While the path is cheap for all solutions, it is expensive when only the  $k$ -sparse solution is desired for a single value  $k$ . The total path (up to  $k$ -sparse solutions) using our computations can be computed in  $O(kmn)$  (assuming finite convergence). While the (approximate) greedy algorithms are more computationally expensive, as will be shown in Section 6, they do offer good practical performance (under measures to be discussed later).

Finally, we again stress the importance of considering such simple approaches where we can only provide convergence to stationary points. Recent convex relaxations for this problem [11, 22], while offering

new insights to the problem, have the same disadvantage that the gap to the optimal solution of problem (22) cannot be computed (together, these methods can give primal-dual gaps). Only a gap to the optimal solution of a relaxation (convex upper bound) is computed. Another major disadvantage is that convex relaxations are not amenable to very large data sets as the per-iteration complexity is  $O(n^3)$  and they require far more iterations in practice. Application of our scheme is limited only by storage of the data; only  $nk$  entries of the covariance matrix are needed at each iteration.

## 5.2 $l_1$ -Constrained PCA

In this section, we focus on the  $l_1$ -constrained PCA problem

$$\max \{x^T A x : \|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}, x \in \mathbf{R}^n\}. \quad (24)$$

In [22], we provide a convex relaxation and corresponding algorithm with a per-iteration complexity of  $O(n^3)$ , but here, we are again only interested in much less computationally expensive methods. We present two recent algorithms that have been proposed through different motivations and show that they can be recovered through our framework.

### An Alternating Maximization Scheme

Recently, Witten et al. [36] have considered the sparse Singular Value Decomposition (SVD) problem

$$\max_{x,y} \{x^T B y : \|x\|_2 = 1, \|y\|_2 = 1, \|x\|_1 \leq k_1, \|y\|_1 \leq k_2, x, y \in \mathbf{R}^n\},$$

where  $B \in \mathbf{R}^{m \times n}$  is the data matrix and  $k_1, k_2$  are positive integers. They recognized that the objective is linear in  $x$  with fixed  $y$  (and vice-versa) and that the problems with  $x$  or  $y$  fixed can be easily solved (see Proposition 15). This fact motivates them to propose a simple cheap alternating maximization scheme which in turn can be used to solve  $l_1$ -constrained PCA. However, the work [36] does not recognize this as yet another instance of the conditional gradient algorithm nor does it provide any convergence results. We show below how to simply recover this algorithm under our framework by applying ConGradU directly to problem (24) for which the convergence claims of Theorem 6 hold.

Thanks to the results established in Section 2, we derive it here with the additional regularization term, i.e., with  $q_\sigma(x) = x^T A_\sigma x$ ,  $\sigma \geq 0$  for an arbitrary matrix  $A \in S^n$  (with  $\sigma = 0$  when  $A$  is already given positive semidefinite). Applying ConGradU, the scheme reads:

$$x^{j+1} = \operatorname{argmax} \{\langle A_\sigma x, x^j \rangle : \|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}, x \in \mathbf{R}^n\}. \quad (25)$$

The above iteration maximizes a continuously differentiable convex function over a compact set so Theorem 6 can be applied to show that every limit point of the sequence converges to a stationary point. Furthermore, thanks to Proposition 15, this scheme reduces to the iteration

$$x^{j+1} = \frac{S_{\lambda^j e}(A_\sigma x^j)}{\|S_{\lambda^j e}(A_\sigma x^j)\|_2} \quad (26)$$

where  $\lambda^j$  is determined by solving (20) with  $a := A_\sigma x^j$  (cf. Proposition 15). The most expensive operation at each iteration is computing  $A_\sigma x^j$  where  $A_\sigma \in \mathbf{S}_+^n$  and  $x^j \in \mathbf{R}^n$  so the per-iteration complexity is  $O(n^2)$ .



### The Expectation-Maximization Algorithm

Another recent approach to problem (24) is developed in [30] which they motivate as an Expectation-Maximization (EM) algorithm for sparse PCA. Motivation comes from their derivation of computing principal components using EM for a probabilistic version of PCA, which is actually equivalent to the power method. The authors solve  $l_1$ -constrained PCA, however also want to enforce that  $\|x\|_0 = k$  at each iteration as well. Their algorithm can be written as

$$x^{j+1} = \frac{S_{\lambda^j e}(A_\sigma x^j)}{\|S_{\lambda^j e}(A_\sigma x^j)\|_2}$$

where  $\lambda^j$  is the  $(k+1)$ -largest entry of vector  $|A_\sigma x^j|$ . Note that the iteration form is identical to that above for the alternating maximization scheme, except for the computation of  $\lambda^j$ . Thus, each iteration can be interpreted as solving

$$x^{j+1} = \operatorname{argmax} \{ \langle A_\sigma x, x^j \rangle : \|x\|_2 = 1, \|x\|_1 \leq \sqrt{k^j}, x \in \mathbf{R}^n \}, \quad (27)$$

where  $k^j$  is chosen at each iteration specifically so that  $x^{j+1}$  is  $k$ -sparse, and can easily be seen to be a variant of ConGradU. Enforcing  $\lambda^j$  to be the  $(k+1)$ -largest entry of vector  $|A_\sigma x^j|$  implicitly sets  $k^j$  to a value that achieves  $k$ -sparsity in  $x^{j+1}$ . While this choice of thresholding enforces exactly  $k$  nonzero entries, the iteration becomes heuristic and neither applies to the true  $l_0$  or  $l_1$ -constrained problem. It is cheap, with the major computation being to compute  $A_\sigma x^j$ , and performs well in practice as shown in Section 6. However, unlike our other iterations, there are no convergence results for this heuristic.

### 5.3 $l_0$ -Penalized PCA

We next consider the  $l_0$ -penalized PCA problem

$$\max \{ x^T A x - s \|x\|_0 : \|x\|_2 \leq 1, x \in \mathbf{R}^n \} \quad (28)$$

which has received most of the recent attention in the literature. We describe two recent algorithms to this problem and show again that they are direct applications of ConGradU.

#### Exploiting Positive Semidefiniteness of $A$

The first approach due to [19] assumes that  $A$  is positive semidefinite (i.e.,  $A = B^T B$  with  $B \in \mathbf{R}^{m \times n}$ ) and writes problem (28) as

$$\max \{ \|Bx\|_2^2 - s \|x\|_0 : \|x\|_2 \leq 1, x \in \mathbf{R}^n \}. \quad (29)$$

The objective is neither concave nor convex. First, using the simple fact (consequence of Lemma 9)

$$\|Bx\|_2^2 = \max_{\|z\|_2 \leq 1} \{ \langle z, Bx \rangle^2 \},$$

the problem is equivalent to

$$\max_{\|x\|_2 \leq 1} \max_{\|z\|_2 \leq 1} \{ \langle z, Bx \rangle^2 - s \|x\|_0 \} = \max_{\|z\|_2 \leq 1} \max_{\|x\|_2 \leq 1} \{ \langle B^T z, x \rangle^2 - s \|x\|_0 \}.$$

Thus, we can now apply Proposition 13 to the inner minimization in  $x$  and then get

$$\max_{x \in \mathbf{R}^n} \{ \|Bx\|_2^2 - s \|x\|_0 : \|x\|_2 \leq 1 \} = \max_{z \in \mathbf{R}^m} \left\{ \sum_{i=1}^n [\langle b_i, z \rangle^2 - s]_+ : \|z\|_2 \leq 1 \right\} \quad (30)$$

where  $b_i \in \mathbf{R}^m$  is the  $i^{th}$  column of  $B$ . This reformulation was previously derived in [10], where the authors also provided a convex relaxation. Note that the reformulation operates in the space  $\mathbf{R}^m$  rather than  $\mathbf{R}^n$ . Since the objective function  $f(z) := \sum_i [\langle b_i, z \rangle^2 - s]_+$  is now clearly convex, we can apply ConGradU. Noting that a subgradient of  $f(z)$  is given by

$$2 \sum_{i=1}^n [\text{sgn}(\langle b_i, z \rangle^2 - s)]_+ (\langle b_i, z \rangle) b_i,$$

the resulting iteration (using Lemma 9) yields:

$$z^{j+1} = \frac{\sum_i [\text{sgn}(\langle b_i, z^j \rangle^2 - s)]_+ (\langle b_i, z^j \rangle) b_i}{\|\sum_i [\text{sgn}(\langle b_i, z^j \rangle^2 - s)]_+ (\langle b_i, z^j \rangle) b_i\|_2}, \quad (31)$$

and the convergence results for the nonsmooth case of Theorem 6 apply.

This is exactly the algorithm recently derived in [19]. Note that an  $O(mn)$  transformation is then needed via Proposition 13 to recover the solution  $x$  of the original problem (29). This is the first cheap ( $O(mn)$  per-iteration complexity) and nongreedy approach for directly solving the  $l_0$ -penalized problem. As with [36] for  $l_1$ -constrained PCA, [29] approach  $l_0$ -penalized PCA via  $l_0$ -penalized SVD. After modifications to write it out for  $l_0$ -penalized PCA, the resulting iteration of their paper is equivalent to iteration (31). However, they did not offer a derivation or state the convergence properties given in [19].

### Approximating the $l_0$ -penalized Problem

As explained in Section 2, we consider the  $l_0$ -penalized problem whereby we use an approximation to the  $l_0$  norm, that is, we consider the problem of maximizing a convex function:

$$\max_x \{x^T A_\sigma x + g(|x|) : \|x\|_2 \leq 1, x \in \mathbf{R}^n\} \quad (32)$$

where the convex function  $g$  is defined by

$$g(z) := -s \sum_{i=1}^n \varphi_p(z_i), \quad s > 0,$$

with  $\varphi_p$  concave satisfying the premises given in Section 2 and  $A_\sigma = A + \sigma I, \sigma \geq 0$ . Applying ConGradU to this problem, as shown in Section 3, using Proposition 8 reduces the iteration to the following weighted  $l_1$ -norm maximization:

$$x^{j+1} = \operatorname{argmax} \{ \langle A_\sigma x, x^j \rangle - \sum_i w_i^j |x_i| : \|x\|_2 = 1 \} \quad (33)$$

where  $w_i^j = s\varphi_p'(|x_i^j|)$ .

Proposition 14 shows that this problem can be solved in closed form so that the conditional gradient algorithm becomes

$$x^{j+1} = \frac{S_{w^j}(A_\sigma x^j)}{\|S_{w^j}(A_\sigma x^j)\|_2}$$

where again  $w_i^j = s\varphi_p'(|x_i^j|)$ . Theorem 6 regarding convergence of every limit point of the resulting sequence to a stationary point again applies. The per-iteration complexity of this iteration is  $O(n^2)$ , which is reduced to  $O(mn)$  if we have the factorization  $A_\sigma = B^T B$ .

Depending on the choice of  $\varphi$  (cf. Section 2), we thus have a family of algorithms for solving problem (28). For example, with the Example 2 (b) given in Section 2, we obtain the recent algorithm of [31] which was derived there by applying what is called the minorization-maximization method, a seemingly different approach; they consider  $A \in \mathbf{S}^n$  and represent the objective as a difference of convex functions plus the penalization:  $x^T Ax - s \sum_i \varphi_p(|x_i|) = x^T A_\sigma x - \sigma x^T x - s \sum_i \varphi_p(|x_i|)$ . The objective is minorized by linearizing the convex term  $x^T A_\sigma x - s \sum_i \varphi_p(|x_i|)$  resulting in a concave lower bound that is maximized. When  $\sigma = 0$ , this is identical to using the conditional gradient algorithm ConGradU, which is only one example of a minorization-maximization method. For more on the minorization-maximization technique and its connection to gradient methods see the recent work [3] and references therein. We also note that [31] derived their algorithm for the sparse generalized eigenvalue (GEV) problem

$$\max \{x^T Ax : x^T Bx \leq 1, \|x\|_0 \leq k, x \in \mathbf{R}^n\} \quad (34)$$

where  $A \in \mathbf{S}^n$  and  $B \in \mathbf{S}_{++}^n$ , which includes as a special case the sparse PCA problem when  $B$  is the identity matrix. The resulting algorithm of [31] for this problem requires computing a matrix pseudoinverse, and is much more computationally expensive (and not amenable to extremely large data sets) than the same algorithm for sparse PCA. Moreover, using the results of Section 2, clearly the general iteration for indefinite  $A$  need not be considered and sparse GEV can always be approached with a closed-form conditional gradient algorithm which still requires computing a matrix pseudoinverse (the closed-form iteration is derived in [31]).

## 5.4 $l_1$ -Penalized PCA

Consider the  $l_1$ -penalized PCA problem

$$\max \{x^T Ax - s\|x\|_1 : \|x\|_2 = 1, x \in \mathbf{R}^n\}. \quad (35)$$

This problem has a nonconvex objective. The work [11] provides a convex relaxation that is solved via semidefinite programming with a per-iteration complexity of  $O(n^3)$ , but here, we are again only interested in much less computationally expensive methods. We describe two methods that exploit positive semidefiniteness of  $A$ , along with a novel scheme that does not require  $A \in \mathbf{S}_n^+$ .

### Reformulation with a Convex Objective

To apply ConGradU, we need either a convex objective or, as shown, an objective of the form  $f(x) + g(\|x\|)$  with  $f, g$  satisfying certain properties (cf. Section 2). Exploiting the fact that  $A \in \mathbf{S}_+^n$ , with  $A = B^T B$ ,  $B \in \mathbf{R}^{m \times n}$ , an equivalent reformulation of the original  $l_0$ -constrained PCA problem can use the square root objective  $\|Bx\|_2$ , and hence the corresponding  $l_1$ -penalized PCA problem reads, instead of (35), as

$$\max \{\|Bx\|_2 - s\|x\|_1 : \|x\|_2 \leq 1, x \in \mathbf{R}^n\}.$$

One can think of this as replacing the objective in our original  $l_0$ -constrained PCA problem (2) with  $\sqrt{x^T Ax} = \|Bx\|_2$  (which is an equivalent problem) and then making the modifications for  $l_1$ -penalized PCA. The objective remains problematic and is neither convex nor concave. However, using again the fact that  $\|Bx\|_2 = \max\{\langle z, Bx \rangle : \|z\|_2 \leq 1, z \in \mathbf{R}^m\}$ , the problem reads

$$\max\{\|Bx\|_2 - s\|x\|_1 : \|x\|_2 \leq 1, x \in \mathbf{R}^n\} = \max_{\|z\|_2 \leq 1} \max_{\|x\|_2 \leq 1} \{\langle z, Bx \rangle - s\|x\|_1\}.$$

Thus, applying Proposition 14, the inner maximization with respect to  $x$  can be solved explicitly and the problem can be reformulated as maximizing a convex objective, and we obtain:

$$\max_{x \in \mathbf{R}^n} \{\|Bx\|_2 - s\|x\|_1 : \|x\|_2 \leq 1\} = \max_{z \in \mathbf{R}^m} \left\{ \sum_{i=1}^n (|b_i^T z| - s)_+^2 : \|z\|_2 \leq 1 \right\},$$

where  $b_i$  is the  $i^{th}$  column of  $B$ .

We can now apply ConGradU to the convex (for similar reasons as for the  $l_0$ -penalized case in Section 5.3) objective  $f(z) = \sum_i [|b_i^T z| - s]_+^2$ , and for which our convergence results for the nonsmooth case hold true. A subgradient of  $f$  is given by

$$2 \sum_{i=1}^n (|b_i^T z| - s)_+ \text{sgn}(b_i^T z) b_i,$$

and, using Lemma 9, the resulting iteration is

$$z^{j+1} = \frac{\sum_i (|b_i^T z^j| - s)_+ \text{sgn}(b_i^T z^j) b_i}{\left\| \sum_i (|b_i^T z^j| - s)_+ \text{sgn}(b_i^T z^j) b_i \right\|_2}. \quad (36)$$

This is exactly the other algorithm recently derived in [19]. Note that to recover the solution to problem (14) needs an  $O(mn)$  transformation via Proposition 14. This algorithm has an  $O(mn)$  per-iteration complexity. Note also that this algorithm was stated earlier in [29] but no such derivation or convergence results were given.

### A Novel Direct Approach

We next derive a novel algorithm for problem (35) by directly applying the conditional gradient algorithm. Indeed, problem (35) reads as maximizing  $f(x) + g(|x|)$  with  $f(x)$  convex,  $g(x)$  convex, differentiable, and monotone decreasing, with  $f(x) = x^T A_\sigma x$  and  $g(u) = -\sum_i u_i$  where  $A_\sigma$  is as previously defined. Applying the ConGradU algorithm and Proposition 8 leads to the iteration

$$x^{j+1} = \operatorname{argmax} \{ \langle A_\sigma x^j, x \rangle - s\|x\|_1 : \|x\|_2 = 1 \}, \quad (37)$$

which by Proposition 14 reduces to

$$x^{j+1} = \frac{S_{se}(A_\sigma x^j)}{\|S_{se}(A_\sigma x^j)\|_2}, \quad (38)$$

where  $e$  is a vector of ones. Theorem 6 applies, showing that any limit point of this iteration is a stationary point of the  $l_1$ -penalized PCA problem.

The matrix-vector product  $Ax^j$  is the main computational cost so the per-iteration complexity is  $O(n^2)$  (or  $O(mn)$  if computing  $B^T(Bx^j)$ ). This approach can handle matrices  $A$  that are not positive semidefinite (by taking  $\sigma > 0$ ) and has stronger convergence results than the conditional gradient method applied to the reformulation of [19], i.e., this approach is equivalent to applying ConGradU with a differentiable objective function (by Proposition 14) and thus satisfies part (c) of Theorem 6. [19] apply ConGradU to a different nondifferentiable formulation for which our theory does not apply.

For the sake of completeness, we end this section by mentioning one of the earlier cheap schemes for sparse PCA, even though it does not fall into the category of directly applying ConGradU.

### An Alternating Minimization Scheme

One of the earlier cheap approaches to sparse PCA, specifically for  $l_1$ -penalized PCA, is proposed in [41] (SPCA). While they generalize all results to multiple factors, we only discuss the one factor case. They pose sparse PCA as an  $l_1/l_2$ -regularized regression problem, specifically

$$(x^*, y^*) = \underset{x, y}{\operatorname{argmin}} \left\{ \sum_{i=1}^m \|b_i - xy^T b_i\|_2^2 + \lambda \|y\|_2^2 + s \|y\|_1 : \|x\|_2^2 = 1, x, y \in \mathbf{R}^n \right\} \quad (39)$$

where  $\lambda$  and  $s$  are the  $l_2$  and  $l_1$  regularization parameters, respectively, and  $B \in \mathbf{R}^{m \times n}$  is a data matrix with rows  $b_i \in \mathbf{R}^n$ . When  $s = 0$ , they show that  $y^*$  is proportional to the leading eigenvector of  $B^T B$ . Indeed, when  $s = 0$ , problem (39) can be recast as a classical maximum eigenvalue problem in  $x$ :

$$x^* = \underset{x}{\operatorname{argmax}} \{x^T B^T B (B^T B + \lambda I_n)^{-1} B^T B x : \|x\|_2^2 = 1, x \in \mathbf{R}^n\} \quad (40)$$

by first solving for  $y$  (simple algebra shows  $y^* = (B^T B + \lambda I_n)^{-1} B^T B x$ ) and plugging  $y^*$  into (39). It is easy to show that the  $x^*$  that solves problem (40) is equal to the leading eigenvector of  $B^T B$  for all  $\lambda \geq 0$ , and thus, for the purposes of finding the leading eigenvector, we do not need to regularize the matrix (i.e., set  $\lambda = 0$ ).

Problem (39) uses an  $l_1$  penalty, known as a LASSO penalty, in order to induce sparsity on  $y$  resulting in an approximate sparse leading eigenvector  $y^*/\|y^*\|_2$ . An alternating minimization scheme in  $x$  and  $y$  is proposed to solve problem (39). For fixed  $y$ , we have

$$\sum_i \|b_i - xy^T b_i\|_2^2 = \sum_i (b_i^T b_i - 2(y^T b_i) b_i^T x + (y^T b_i)^2 x^T x) = -2y^T \left( \sum_i b_i b_i^T \right) x + C$$

where  $C$  is a constant (using the constraint  $\|x\|_2 = 1$ ), so that the minimizer  $x^*$  is solved by maximizing a linear function over the unit sphere which, by Lemma 9, is easily solved in closed-form. For fixed  $x$ , the minimizer  $y^*$  is found by solving an unconstrained minimization problem of the form  $\|\cdot\|_2^2 + s\|\cdot\|_1$  (also known as the *elastic net* problem). This problem can be solved efficiently for fixed  $s$  using fast first-order methods such as FISTA [2] or for a full path of values for  $s$  using LARS [14]. Thus, [41] solve a nonconvex problem in two variables using alternating minimization. While this scheme is computationally inexpensive compared to convex relaxations, it is not as cheap as the schemes we are considering due to the subproblem with fixed  $x$ , and no convergence results have been derived for it.

## 6 Experiments

Thus far, various algorithms have been provided with the goal of learning sparse rank one approximations. In this section, these different methods are compared. The algorithms considered here are  $l_0$ -constrained PCA (novel iteration), an approximate greedy algorithm [10], GPowerL1 ( $l_1$ -penalized PCA of [19]), GPowerL0 ( $l_0$ -penalized PCA of [19]), Expectation-Maximization ( $l_1$ -constrained PCA of [30]), and thresholding (select  $k$  entries of principal eigenvector with largest magnitudes). We also consider an exact greedy algorithm and the optimal solution (via exhaustive search) for small dimensions ( $n = 10$ ).

The goal of these experiments is two-fold. Firstly, we demonstrate that the various algorithms give very similar performance. The measure of comparison used is the proportion of variance explained by a sparse vector versus that explained by the true principal eigenvector, i.e., the ratio  $x^T A x / v^T A v$  where  $x$  is the sparse eigenvector and  $v$  is the true principal eigenvector of  $A$ . The second goal is to solve very large

sparse PCA problems. The largest dimension we approach is  $n = 50000$ , however, as discussed above, the ConGradU algorithm applied to  $l_0$ -constrained PCA has very cheap  $O(mn)$  iterations and is limited only by storage of a data matrix. Thus, on larger computers, extremely large-scale sparse PCA problems (much larger than those solved even here) are also feasible.

Note that we do not compare against all algorithms listed in Table 1. In particular, SPCA [41] was already demonstrated to be computationally more expensive, as well to provide inferior performance, to GPowerL1 and GPowerL0 [19]. The  $l_1$ -constrained PCA method of [36] and  $l_0$ -penalized PCA method of [31] are also cheap methods that give similar performance (learned from experiments not shown in this paper) to the algorithms in our experiments. Finally, note that, for all experiments, we do a postprocessing step in which we compute the largest eigenvector of the data matrix in the  $k$ -dimensional subspace that is discovered by the respective methods.

All experiments were performed in MATLAB on a PC with 2.40GHz processor with 3GB RAM. Codes from the competing methods were downloaded from URL's available in the corresponding references. Slight modifications were made to do singular value decompositions rather than eigenvalue decompositions in order to deal with much smaller  $m \times n$  data matrices rather than  $n \times n$  covariance matrices. We first demonstrate performance on random matrices and follow with a text data example.

## 6.1 Random Data

We here consider random data matrices  $F \in \mathbf{R}^{m \times n}$  with  $F_{ij} \sim N(0, 1/m)$ . It was already shown in literature on greedy methods [10] and convex relaxations [11, 22] that random matrices of the form  $xx^T + U$  where  $U$  is uniformly distributed noise are *easy* examples. Results here show that taking sparse eigenvectors of the matrix  $F^T F$  is also relatively easy.

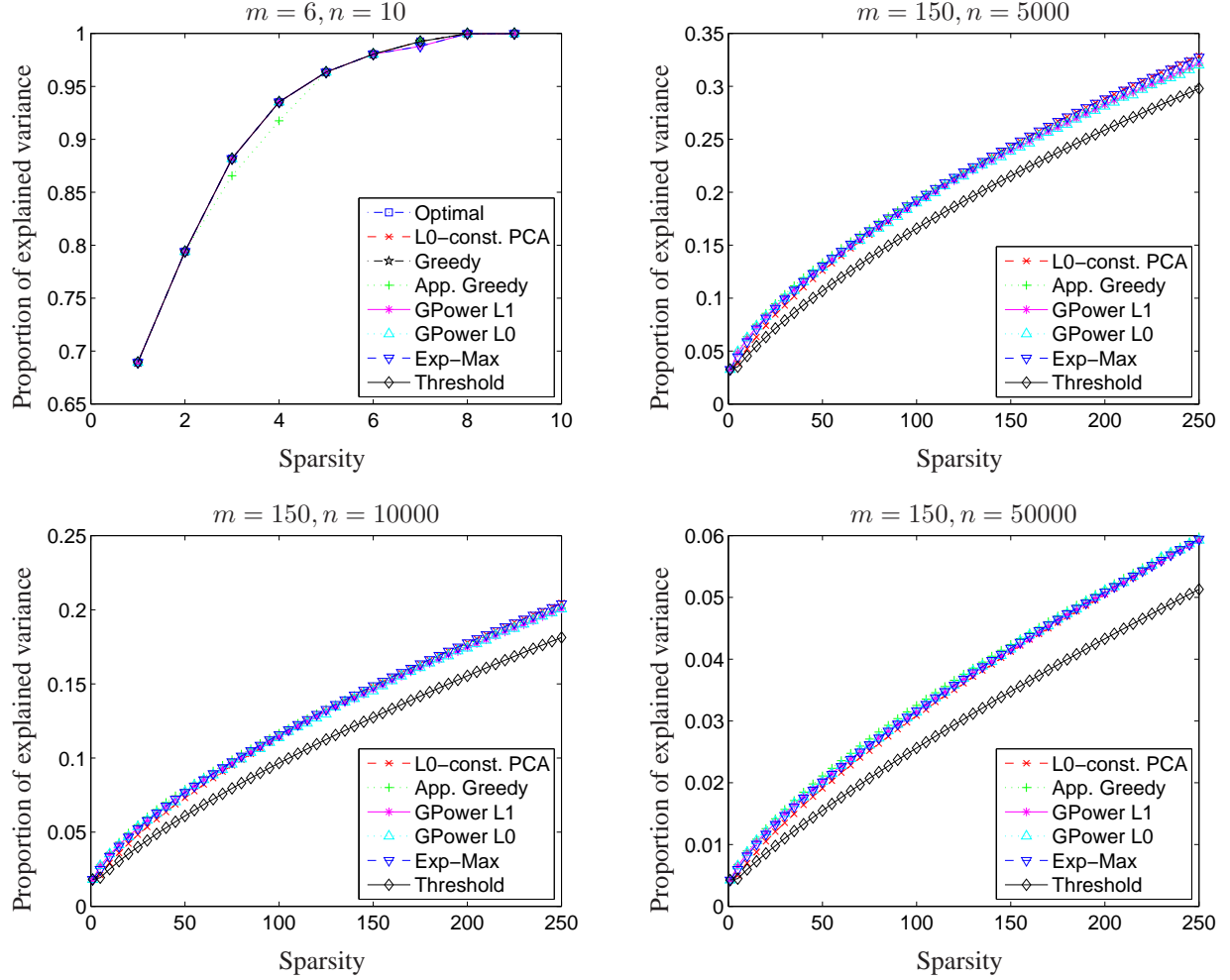
The experiments consider  $n = 10$  ( $m = 6$ ) and  $n = 5000, 10000, 50000$  (each with  $m = 150$ ), each using 100 simulations. We consider  $l_0$ -constrained PCA with  $k = 2, \dots, 9$  for  $n = 10$  and  $k = 5, 10, \dots, 250$  for the remaining tests. The optimal solution (found by exhaustive search) and the exact greedy algorithm (too computationally expensive for high dimensions) are only used when  $n = 10$ .

Figure 2 compares performance of the various algorithms. Similar patterns are seen as  $n$  increases. For  $n = 10$ , optimal performance is obtained for almost every algorithm. For higher dimensions, we have no measure of the gap to optimality. As the dimension increases, the proportion of explained variation by using the same fixed cardinality decreases as expected. The next subsection shows that we do not necessarily need to explain most of the variation in the true eigenvector in order to gain interpretable factors.

Results only up to a cardinality level of 250 variables are displayed because our goal is simply to compare the different algorithms. All algorithms, except for simple thresholding, perform very similarly. These figures do not sufficiently display the story, so we describe the similar pattern that occurs. The approximate greedy algorithm does best for smallest cardinalities, then the expectation-maximization scheme dominates, and at some point the novel  $l_0$ -constrained PCA scheme gives best performance at a higher level of explained variation. For  $n = 50000$ , we do not actually see this change yet because such little variation is explained with only 250 variables. Furthermore, it is important to notice that the thresholded solution is consistently and greatly outperformed by all other methods, suggesting that the performance results are enhanced via the conditional gradient algorithm. These experiments simply show that these algorithms offer very similar performance, and hence we next compare them computationally.

Figure 3 displays the computational time comparing the various algorithms for the different dimensions. Firstly, note that for penalized PCA problems (GPowerL0 and GPowerL1), a parameter must be tuned in order to achieve the desired sparsity. Figures here do not account for time spent tuning parameters. Secondly, the greedy algorithms must compute the greedy solution at all sparsity levels of  $1, \dots, 250$  in

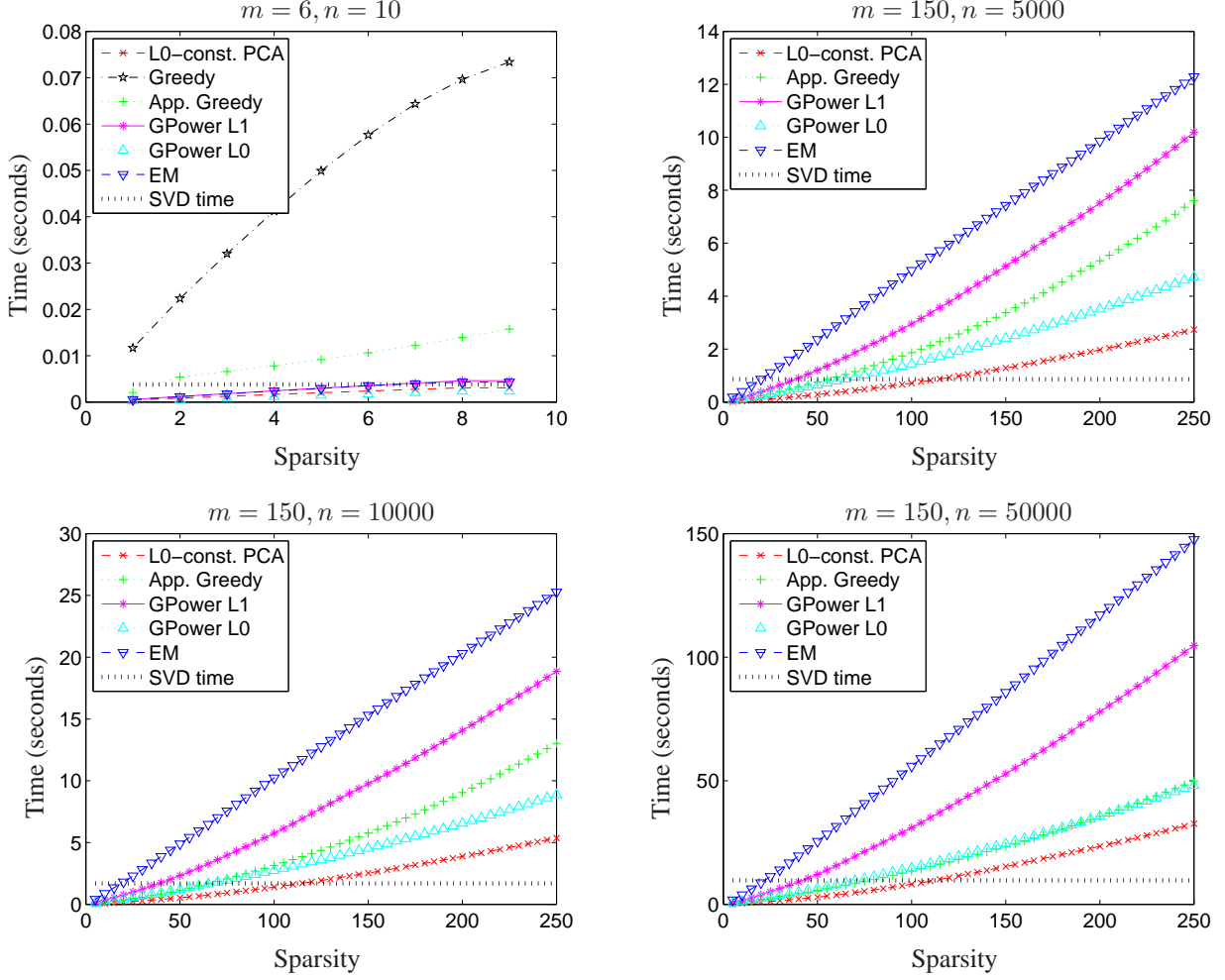




**Figure 2:** Plots show the average percent of variance explained by the sparse eigenvectors found by several algorithms, i.e., the ratio  $x^T(F^T F)x/v^T(F^T F)v$  where  $x$  is the sparse eigenvector and  $v$  is the true first eigenvector.  $F$  is an  $m \times n$  matrix with  $F_{ij} \sim N(0, 1/m)$ . Sparsities of  $2, \dots, 9$  are computed for  $n = 10$  and  $5, 10, \dots, 250$  for the remaining experiments. 100 simulations are used to produce all results.

order to obtain the solution with 250 variables, and hence the time displayed to compute greedy solutions is the cumulative time. Time for each algorithm is thus also taken as the cumulative time, albeit for all others it is the cumulative time to obtain a solution with sparsity levels of  $5, 10, \dots, 250$ . The novel  $l_0$ -constrained algorithm requires an initial solution which we take as thresholded solution of the true principal eigenvector. The time to obtain that initial solution is marked as `svdTime`, however it need only be computed once for the entire path.

For  $n = 10$ , the exact greedy algorithm is clearly the most expensive, requiring  $n$  maximum eigenvalue computations per iteration. For higher dimensions, the same pattern occurs. The expectation-maximization scheme requires the most time because the scheme implicitly solves a penalized problem and thus also implicitly tunes a parameter. GPowerL1 is surprisingly (since the tuning time is not included) next. Despite



**Figure 3:** Plots show the average cumulative time to produce the sparse eigenvectors of  $F^T F$  found by several algorithms, with  $F$  an  $m \times n$  matrix with  $F_{ij} \sim N(0, 1/m)$ . Sparsities of  $2, \dots, 9$  are computed for  $n = 10$  and  $5, 10, \dots, 250$  for the remaining experiments. Note that cumulative times are given, i.e., the time to calculate the vector with 30 nonzeros adds up the time to compute vectors with  $5, 10, \dots, 30$  nonzeros, in order to compare with the approximate greedy method. The  $\text{svdTime}$  is the time required to compute the principal eigenvector of  $F^T F$  which is used to compute an initial solution for  $l_0$ -constrained PCA. 100 simulations are used to produce all results.

being cheap, it requires more iterations than other methods to converge. The approximate greedy algorithm follows, and is expected to be (relatively) computationally expensive because of the maximum eigenvalue computed at each iteration. This is followed by GPowerL0 and finally by the cheapest scheme, the novel  $l_0$ -constrained PCA iteration.

We now discuss the advantages and disadvantages of the different schemes. Clearly, if the sparsity is known (or the sparsities desired is much less than the full path), the approximate greedy algorithm is much more computationally expensive. Comparing the other cheap schemes, the  $l_0$ -constrained PCA scheme is cheapest (given the initial solution). The disadvantage of the penalized schemes (GPowerL1 and GPowerL0)

is that they must be tuned which is computationally very expensive (not shown). Warmstarting could be used, for example, by initializing for  $k = 10$  based on the solution to  $k = 5$  rather than from the thresholded solution. Thus, if the desired sparsity is known, the  $l_0$ -constrained PCA scheme is clearly the algorithm to use. If not, then all of the algorithms are cheap, offer similar performance, and can be used to derive a path of sparse solutions.

## 6.2 Republicans or Democrats: What is the Difference?

We consider here text data based on all State of the Union addresses from 1790-2011. Transcripts are available at <http://stateoftheunion.onetwothree.net> where other interesting analyses of this data are also done. Here, sparse PCA is used to further analyze these historical speeches. Questions one might ask relate to how the language in speeches has changed from George Washington through Barak Obama or how the relevant issues divide the different presidents. After taking the stems of words and removing commonly used *stopwords*, we created a bag-of-words data set based on all remaining words, leaving 12953 words (i.e.,  $n = 12953$  for this example). Our data matrix here is  $B \in \mathbf{R}^{m \times n}$  where  $m$  is the number of speeches and  $B_{ij}$  is the number of times the  $j^{\text{th}}$  word occurs in the  $i^{\text{th}}$  State of the Union address. We analyze two different sample sizes: using all speeches from 1790-2011 ( $m = 225$ ) and only speeches from 1982-2011 ( $m = 31$ ). The rows of  $B$  are normalized so that each speech is of the same length. The following results are for PCA and sparse PCA performed on the covariance matrix  $A = B^T B$ .

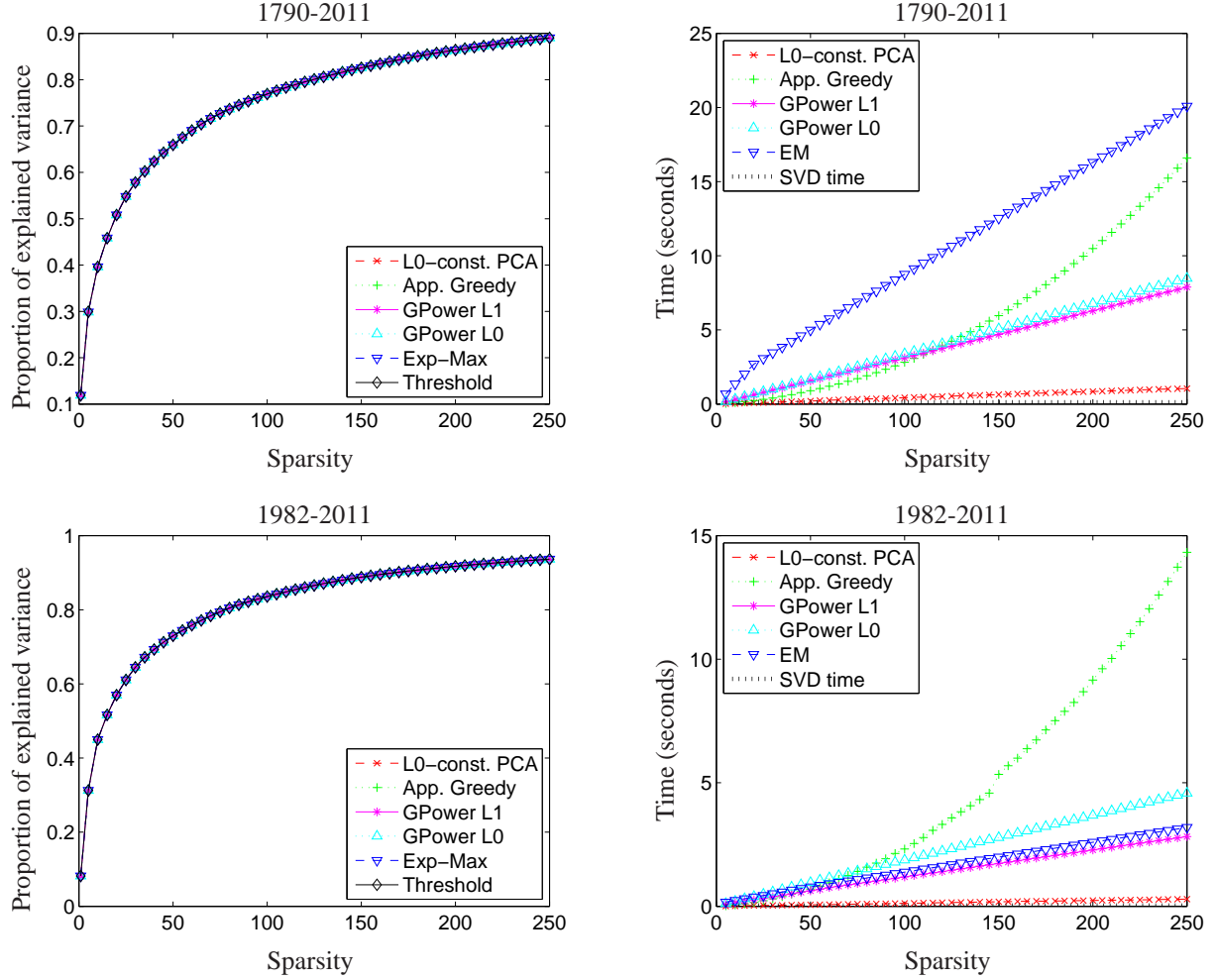
Figure 4 displays the performance of the various algorithms on the text data set using the same measure as with random data above. For this high-dimensional data set, much fewer variables are needed to explain more variation relative to what was observed with the random matrices. Another major difference is that each algorithm gives exactly the same solution; even just taking the thresholded solution gives the same solution. This leads to postulate that real data often might contain a structure that makes it rather easy to solve (still, of course, we have no results on solution quality). Note that the  $l_0$ -constrained scheme seems to be the cheapest algorithm here because, starting at the thresholded solution, it required one iteration to achieve convergence! Despite the simplicity of obtaining sparse solutions for this data, we continue to show what can be learned using this tool. The goal is to show that sparse factors offer interpretability that cannot be learned from using all 12953 variables.

Figure 5 (left) shows the result of projecting the data on the first two principal eigenvectors<sup>8</sup>. The second factor clearly clusters the speeches into two groups (roughly into positive versus negative coordinates in the second factor). Examination of these two groups shows a chronological pattern which as seen in the figure clusters those speeches that occurred before (and during) World War I with those speeches that were made after the war. Figure 5 (right) shows the data projected on sparse factors giving very a similar illustration. Factors 1 and 2 use 150 and 15 variables, respectively. Table 2 displays the words from the sparse second factor and their sign. We *roughly* associate the positively weighted words with speeches before the war and negatively weighted words with speeches after the war. Then one might interpret that, before World War I, presidents spoke about the United States of America as a collection of united states, but afterwards, spoke about the country as one american nation that faced issues as a whole.

Figure 6 next shows a similar analysis using only speeches from 1982-2011. A clear distinction between republicans and democrats was discovered. Again, sparse PCA is used to interpret the factors with 15 variables each. Table 3 shows the most important words discovered by PCA (using thresholding) and sparse PCA for the top 3 factors. The first factor gives the same words for both analyses and no clear interpretation is seen. The second factors are more interesting. The thresholded PCA factor clearly relates to international

---

<sup>8</sup>We use projection deflation, as described in [23], to obtain multiple sparse factors.

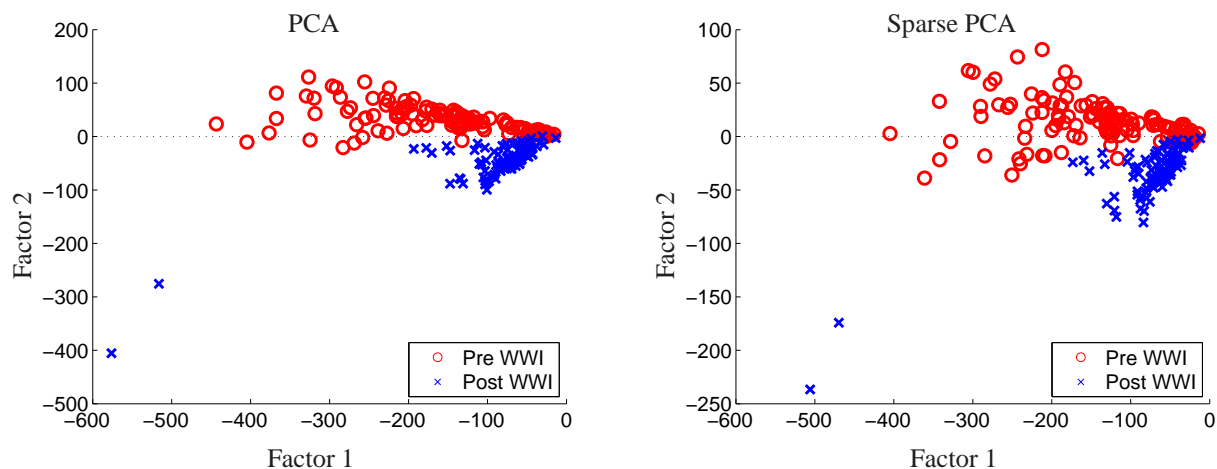


**Figure 4:** The left plots shows the average percent of variance explained by the sparse eigenvectors and the right plots show the computational time to run the algorithms. Data is from the text of State of the Union addresses. Two different numbers of samples are used: all addresses from 1790-2011 and just those from 1982-2011.

security issues. The sparse PCA factor, however, clearly focuses on domestic issues, e.g., health-care and education. Differences occur because PCA deflates with the true eigenvectors while sparse PCA deflates with the sparse factors. In any case, these are clearly issues that divide republicans and democrats. The third thresholded PCA factor seems to be related to educational reforms and funding them. The third sparse PCA factor is clearly related to fiscal policies and the economy.

## 7 Concluding Remarks

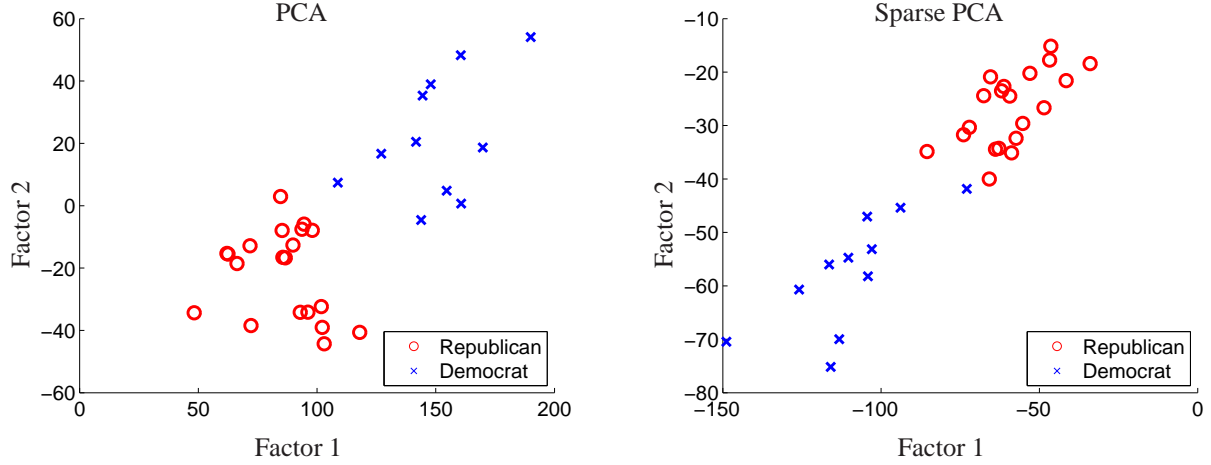
Sparse PCA admits a simple formulation which maximizes a (usually convex) quadratic objective subject to *seemingly* simple constraints. Nonetheless, it is a difficult nonconvex optimization problem, and a large literature has focused on various modifications of the desired problem in order to derive simple algorithms



**Figure 5:** The left plot shows the text of all State of the Union addresses from 1790-2011 reduced from 12953 to 2 dimensions using PCA. The right plot shows the same data reduced to 2 dimensions using  $l_0$ -constrained PCA. Factor 1 is a function of 150 words, while Factor 2 is a function of only 15 words.

Factor 2	Sign
govern	+
state	+
unite	+
american	-
econom	-
feder	-
help	-
million	-
more	-
nation	-
new	-
program	-
work	-
world	-
year	-

**Table 2:** The table shows the top 15 out of 12953 words associated with the second factor derived from  $l_0$ -constrained PCA on the text of all State of the Union addresses from 1790-2011. The words for Factor 2 of thresholded PCA are almost identical. A first factor with 150 nonzero entries was deflated from the data. The signs of the words are given to show what drives the difference between the clusters in Figure 5. Note that words here are *word stems*, i.e. the words *program* and its plural *programs* are both represented by *program*.



**Figure 6:** The left plot shows the text of 31 State of the Union addresses from 1982-2011 reduced from 12953 to 2 dimensions using PCA. The right plot shows the same data reduced to 2 dimensions using  $l_0$ -constrained PCA where both factors are functions of exactly 15 words.

PCA Factors 1-3			Sparse PCA Factors 1-3		
1	2	3	1	2	3
year	america	govern	year	job	budget
american	work	peopl	american	countri	economi
more	terrorist	program	more	children	program
peopl	world	school	peopl	secur	live
work	freedom	centuri	work	famili	million
america	nation	commun	new	tonight	over
new	peopl	spend	america	care	busi
nation	more	children	nation	health	futur
make	cut	america	make	ask	plan
help	iraq	new	help	last	reform
govern	year	tax	govern	school	most
world	terror	countri	world	state	mani
time	care	deficit	time	support	respons
tax	job	challeng	congress	commun	good
congress	secur	support	tax	cut	invest

**Table 3:** The left side of the table shows the top 15 words associated with the first 3 factors derived from PCA on the text of 31 State of the Union addresses from 1982-2011. The right side of the table shows the top 15 words when the 3 factors are derived using  $l_0$ -constrained PCA. The original number of dimensions is 12953. Note that words here are *word stems*, i.e. the words *program* and its plural *programs* are both represented by *program*.

that hopefully produce *good* approximate solutions. No gaps to the solution of the original problem are



given by any of the currently known schemes.

In this paper, we have shown that the conditional gradient algorithm ConGradU:

- can be directly applied to the original  $l_0$ -constrained PCA problem (2) without any modifications to produce a very simple scheme with low computational complexity.
- can be successfully applied to maximizing a convex function over an arbitrary compact set and was proven to exhibit global convergence to stationary points. Efficiency of this scheme builds on the result that, while maximizing a quadratic function over the  $l_2$  unit ball with an  $l_0$  constraint is a difficult problem, maximizing a linear function over the same nonconvex set is simple.
- provides a unifying framework to derive and analyze all new and old schemes discussed in the paper which, as we have seen, were derived from disparate approaches in the cited literature. As shown, all these algorithms which have been used in various applications are special cases of ConGradU.

The overall message is that, for some difficult problems, we can achieve the same (limited) theoretical guarantees and practical performance using the same algorithm on modified (seemingly easier) problems as we can on the original difficult problem. Furthermore, all of these algorithms emerging from ConGradU give similar performance in practice and with similar complexities.

We conclude by showing that the same tools we have used for applying the ConGradU algorithm to sparse PCA can readily be applied to other sparsity-constrained statistical tools, e.g., sparse Singular Value Decompositions, sparse Canonical Correlation Analysis, and sparse nonnegative PCA. Note that our comments below about  $l_0$ -constrained problems can also be extended to the corresponding  $l_1$ -constrained and  $l_0/l_1$ -penalized problems.

### **Sparse Singular Value Decomposition (SVD)**

Sparse SVD solves the problem

$$\max \{x^T B y : \|x\|_2 = 1, \|y\|_2 = 1, \|x\|_0 \leq k_1, \|y\|_0 \leq k_2, x \in \mathbf{R}^m, y \in \mathbf{R}^n\} \quad (41)$$

where  $B \in \mathbf{R}^{m \times n}$  is the data matrix and  $k_1, k_2$  are positive integers. Note that this is equivalent to  $l_0$ -constrained PCA when  $x = y$  and  $B \in \mathbf{S}^n$ . [36] considered a relaxation to this problem as described in Section 5.2. They relaxed the  $l_0$  constraints on  $x$  and  $y$  with  $l_1$  constraints (and relaxed equalities to inequalities) and applied an alternating maximization scheme, where each optimization problem is easily solved via Proposition 15. Following exactly the approach we suggested for  $l_0$ -constrained PCA, there is no need to relax the  $l_0$  constraints. Proposition 11 can be used to solve directly the alternating optimization problems, giving rise to a simple algorithm for the  $l_0$ -constrained SVD problem (41).

### **Sparse Canonical Correlation Analysis (CCA)**

Sparse CCA solves the problem

$$\max \{x^T B^T C y : x^T B^T B x = 1, y^T C^T C y = 1, \|x\|_0 \leq k_1, \|y\|_0 \leq k_2, x \in \mathbf{R}^p, y \in \mathbf{R}^q\} \quad (42)$$

where  $B \in \mathbf{R}^{m \times p}$  and  $C \in \mathbf{R}^{m \times q}$  are data matrices and  $k_1, k_2$  are positive integers. [36] suggest that useful (i.e., interpretable) results can still be obtained by substituting the identity matrix for  $B^T B$  and  $C^T C$  in the constraints, resulting in the sparse SVD problem above. Rather, we propose substituting the diagonals of  $B^T B$  and  $C^T C$  as proxies. Propositions 11 and 15 (among others in Section 4) can both easily be extended to optimizing over the constraints  $x^T D x = 1$  and  $x^T D x \leq 1$ , where  $D$  is diagonal. A simple

alternating maximization algorithm then follows for the resulting  $l_0$ -constrained *approximate* CCA problem.

### Sparse Nonnegative Principal Component Analysis (PCA)

Sparse nonnegative PCA solves the problem

$$\max \{x^T A x : \|x\|_2 = 1, \|x\|_0 \leq k, x \geq 0, x \in \mathbf{R}^n\} \quad (43)$$

where  $A \in \mathbf{R}^{n \times n}$  is a covariance matrix and  $k$  is a positive integer. This is exactly the  $l_0$ -constrained PCA problem (2) with additional nonnegativity constraints. Simple extensions of Propositions 11 and 15 lead to a simple scheme for  $l_0$ -constrained nonnegative PCA based on the ConGradU algorithm.

## 8 Acknowledgements

The authors thank Nouredine El Karoui for useful discussions regarding random matrices. This research was partially supported by the United States-Israel Science Foundation under BSF Grant #2008-100.

## References

- [1] O. Alter, P. Brown, and D. Botstein. Singular value decomposition for genome-wide expression data processing and modeling. *Proceedings of the National Academy of Sciences*, 97:10101–10106, 2000.
- [2] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal of Imaging Sciences*, 2(1):183–202, 2009.
- [3] A. Beck and M. Teboulle. Gradient-based algorithms with applications to signal recovery problems. *Convex Optimization in Signal Processing and Communications*, pages 42–88, 2010. Edited by Yonina Eldar and Daniel Palomar.
- [4] D. Bertsekas. *Nonlinear Programming, 2nd Edition*. Athena Scientific, 1999.
- [5] M. Blum, R. W. Floyd, V. Pratt, R. Rivest, and R. Tarjan. Time bounds for selection. *Journal of Computer and System Sciences*, 7(4):448–461, 1973.
- [6] R. N. Bracewell. *The Fourier Transformation and its Applications*. The McGraw Hill Companies, Inc., 3rd edition, 2000.
- [7] J. Cadima and I. T. Jolliffe. Loadings and correlations in the interpretation of principal components. *Journal of Applied Statistics*, 22(2):203–214, 1995.
- [8] E. Candès and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005.
- [9] P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling and Simulation*, 4(4):1168–1200, 2005.
- [10] A. d’Aspremont, F. Bach, and L. El Ghaoui. Optimal solutions for sparse principal component analysis. *Journal of Machine Learning Research*, 9:1269–1294, 2008.
- [11] A. d’Aspremont, L. El Ghaoui, M.I. Jordan, and G. R. G. Lanckriet. A direct formulation for sparse PCA using semidefinite programming. *SIAM Review*, 48(3):434–448, 2007.
- [12] J. C. Dunn. Convergence rates for conditional gradient sequences generated by implicit step length rules. *SIAM Journal of Control and Optimization*, 18(5), 1979.
- [13] J. C. Dunn. Convergence rates for conditional gradient sequences generated by implicit step length rules. *SIAM Journal of Control and Optimization*, 18(5), 1980.

- [14] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32(2):407–499, 2004.
- [15] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3:95–110, 1956.
- [16] P.J. Hancock, A.M. Burton, and V. Bruce. Face processing: human perception and principal components analysis. *Memory and Cognition*, 26, 1996.
- [17] L.J. Hargrove, G. Li, K.B. Englehart, and B.S. Hudgins. Principal components analysis preprocessing for improved classification accuracies in pattern-recognition-based myoelectric control. *IEEE Transactions on Biomedical Engineering*, 56(5):1407–1414, 2009.
- [18] I. T. Jolliffe, N.T. Trendafilov, and M. Uddin. A modified principal component technique based on the LASSO. *Journal of Computational and Graphical Statistics*, 12(3):531–547, 2003.
- [19] M. Journée, Y. Nesterov, P. Richtárik, and R. Sepulchre. Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, 11:517–553, 2010.
- [20] A. Lakhina, M. Crovella, and C. Diot. Diagnosing network-wide traffic anomalies. *Proceedings of ACM Conference of the Special Interest Group on Data Communication (SIGCOMM)*, 2004.
- [21] E. S. Levitin and B. T. Polyak. Constrained minimization methods. *USSR Comp Math and Math Phys*, 6:1–50, 1966.
- [22] R. Luss and M. Teboulle. Convex approximations to sparse pca via lagrangian duality. *Operations Research Letters*, 39(1):57–61, 2011.
- [23] Lester Mackey. Deflation methods for sparse pca. *Advances in Neural Information Processing Systems*, 21:1017–1024, 2009.
- [24] O. L. Mangasarian. Machine learning via polyhedral concave minimization. In H. Fischer, B. Riedmueller, and S. Schaeffer, editors, *Applied Mathematics and Parallel Computing - Festschrift for Klaus Ritter*, Physica-Verlag A Springer-Verlag Company, Heidelberg, pages 175–188, 1996.
- [25] B. Moghaddam, Y. Weiss, and S. Avidan. Spectral bounds for sparse PCA: Exact and greedy algorithms. *Advances in Neural Information Processing Systems*, 18:915–922, 2006.
- [26] J. J. Moreau. Proximité et dualité dans un espace hilbertien. *Bull. Soc. Math. France*, 93:273–299, 1965.
- [27] A.L. Price, N.J. Patterson, R.M. Plenge, M.E. Weinblatt, N.A. Shadick, and D. Reich. Principal components analysis corrects for stratification in genome-wide association studies. *Nature Genetics*, 38(8):904–909, 2006.
- [28] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [29] H. Shen and J. Huang. Sparse principal component analysis via regularized low rank matrix approximation. *Journal of Multivariate Analysis*, 101:1015–1034, 2008.
- [30] C. D. Sigg and J. M. Buhmann. Expectation-maximization for sparse non-negative pca. *Proceedings of the 25th International Conference on Machine Learning*, 2008. 8 pages.
- [31] B. K. Sriperumbudur, D. A. Torres, and G. R. G. Lanckriet. A majorization-minimization approach to the sparse generalized eigenvalue problem. *Machine Learning*, 2010. DOI 10.1007/s10994-010-5226-3.
- [32] G. Strang. *Linear Algebra and Its Applications*, 4th Edition. Brooks Cole, 2005.
- [33] R. Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal statistical society, series B*, 58(1):267–288, 1996.
- [34] N.T. Trendafilov and I. T. Jolliffe. Projected gradient approach to the numerical solution of the SCoTLASS. *Journal of Computational Statistics and Data Analysis*, 50:242–253, 2006.

- [35] J. Weston, A. Elisseeff, B. Schölkopf, and M. Tipping. Use of the zero-norm with linear models and kernel methods. *Journal of Machine Learning Research*, (2):1439–1461, 2003.
- [36] D. M. Witten, R. Tibshirani, and T. Hastie. A penalized matrix decomposition, with applicaitons to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534, 2009.
- [37] W. I. Zangwill. *Nonlinear Programming: A Unified Apporach*. Englewood Cliffs, N. J., Prentice-Hall, 1969.
- [38] J. Zhang, Y. Yan, and M. Lades. Face recognition: Eigenface, elastic matching, and neural nets. *Proceedings of the IEEE*, 85(9), 1997.
- [39] Z. Zhang, H. Zha, and H. Simon. Low-rank approximations with sparse factors i: Basic algorithms and error analysis. *SIAM Journal on Matrix Analysis and Applications*, 23(3):706–727, 2002.
- [40] Z. Zhang, H. Zha, and H. Simon. Low-rank approximations with sparse factors ii: Penalized methods with discrete newton-like iterations. *SIAM Journal on Matrix Analysis and Applications*, 25(4):901–920, 2004.
- [41] H. Zou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15(2):265–286, 2006.